

# Finite sample corrections for parameters estimation and significance testing

Boon Kin Teh<sup>a,b,1</sup>, Darrell JiaJie Tay<sup>a,b</sup>, Sai Ping Li<sup>c</sup>, and Siew Ann Cheong<sup>a,b</sup>

<sup>a</sup>Division of Physics and Applied Physics, School of Physical and Mathematical Sciences, Nanyang Technological University, 21 Nanyang Link, Singapore 637371, Republic of Singapore; <sup>b</sup>Complexity Institute, Block 2 Innovation Centre, Level 2 Unit 245, Nanyang Technological University, 18 Nanyang Drive, Singapore 637723, Republic of Singapore.; <sup>c</sup>Institute of Physics, Academia Sinica, 128 Section 2, Academia Road, Nankang, Taipei 11529, Taiwan.

This manuscript was compiled on May 22, 2017

An increasingly important problem in the era of Big Data is fitting data to distributions. However, many stop at visually inspecting the fits or use the coefficient of determination as a measure of the goodness of fit. In general, goodness-of-fit measures do not allow us to tell which of several distributions fit the data best. Also, the likelihood of drawing the data from a distribution can be low even when the fit is good. To overcome these limitations, Clauset et al. advocated a three-step procedure for fitting any distribution: (i) estimate parameter(s) accurately, (ii) choosing and calculating an appropriate goodness of fit, (iii) test its significance to determine how likely this goodness of fit will appear in samples of the distribution. When we perform this significance testing on exponential distributions, we often obtain low significance values despite the fits being visually good. This led to our realization that most fitting methods do not account for effects due to the finite number of elements and the finite largest element. The former produces sample size dependence in the goodness of fits and the latter introduces a bias in the estimated parameter and the goodness of fit. We propose modifications to account for both and show that these corrections improve the significance of the fits of both real and simulated data. In addition, we used simulations and analytical approximations to verify that convergence rate of the estimated parameters towards its true value depends on how fast the largest element converge to infinity.

Significance Testing | Finite Sample Effects | Curve Fitting | Maximum Likelihood

The current era of Big Data has ushered in a new way to look at Science — and that is letting the data speak for itself. Because of this, we are now much more concerned about empirical distributions than we have in the past, and to check what the empirical distributions could be in statistically rigorous ways. In the past, statisticians commonly tested empirical data against the univariate normal distribution (1). Some of these tests focus on the goodness-of-fit of higher order moments (2–4), while others compare the test statistics against an Empirical Distribution Function (EDF) (5–8). In 2011, Nornadiah and Yap performed a systematic comparison of Anderson-Darling (AD), Lilliefors, Kolmogorov-Smirnov (KS), and Shapiro-Wilk (SW), using numerical simulation and concluded that the SW test is the best, followed closely by the AD test for a given significance (9).

Among these tests, the KS and Lilliefors tests can also be applied to non-normal distributions. In fact, many real-world data do not follow normal distributions. For instance, many social systems are known to have power-law distributions (10). These include the financial returns (11–14), word count (15, 16), city size (17, 18), home price (19–21), wealth and income (22, 23) distributions. One simple but naive way to detect a power law is to plot the data in log-log scale, fit it to a

straight line and determine the goodness of fit. However, this simple method has three major flaws: (i) many distributions (eg. exponential, gamma, log-normal) can also look straight in log-log plot, especially if the range of data is small; (ii) the goodness of fit only quantifies how well the fit is visually but does not tell us how plausible the fit is; and (iii) if our data looks straight in both log-log and semi-log plots, the goodness of fit values obtained from the two cannot be directly compared since they were obtained from plots of different scales. Clauset, Shalizi and Newman (CSN) address precisely these three points in their 2009 paper (24), and the test they proposed is now considered by many the gold standard in curve fitting. We shall describe the main idea of the CSN technique in greater mathematical detail in Section 1.

Since the CSN test can be applied across distributions, we also use it to fit data that appear exponentially distributed. On many occasions, we discovered that the exponential fits look good visually, but have significance values (*p*-value) much lower than fits of other data to power laws, even though the latter look visibly poorer. In fact, in the CSN paper where empirical data is tested against a power law (PL), log-normal, exponential (EXP), stretched exponential, and a power law with cut-off, the exponential distribution consistently performs poorer than the other distributions. This was also the case when Brzezinski tested the upper-tail wealth data for China, Russia, US, and the World using the CSN method (25). In these papers, the data might truly be non-exponentially distributed, so it is not surprising the exponential fits fail. However, the low *p*-values for the visually convincing exponen-

## Significance Statement

All real data sets are finite, hence they suffer from finite sample effects. These effects cannot be neglected, as they can lead us to sometimes underestimate the statistical significance of the fits and reject fits that are visibly good. Here, we propose modifications and outline a procedure for curve fitting that accounts for these effects, allowing us to compare fits across different distributions and different sample sizes. Indeed, after the changes we see significant improvement in the statistical significance of the visually good fits, reducing the risk of mistakenly rejecting plausible fits. This will benefit all disciplines fitting large data sets to several distribution functions, and use the best candidate to guide the development of theories and models.

B.K.T., J.J.D.T., and S.A.C. designed research; B.K.T. performed research; B.K.T., J.J.D.T. and S.P.L. collected data; B.K.T. and J.J.D.T. analyzed data; all authors wrote and reviewed the paper.

The authors declare no conflict of interest.

<sup>1</sup>To whom correspondence should be addressed. E-mail: cheongsa@ntu.edu.sg

tial fits to our data suggest that something fundamental was missed.

We realized there are two issues associated with fitting data to distributions defined over  $(0, \infty)$ . First, there is the *finite largest element effect* (FLE), due to the largest element in the data being finite. Second, we also encounter the *finite number of elements effect* (FNE), due to the sample size dependence of the goodness-of-fit measures. These two *finite sample effects* are often neglected in tests of statistical significance. After describing the CSN test, we illustrate in Section 1 the FLE and FNE effects by applying the test to three real data sets. With the insights gained, we designed both the estimators and test statistic to account for the FLE and FNE effects in Section 2. Since real data is frequently polluted by noise, we also discuss the impact of noise on the  $p$ -value, and propose a test statistic that accounts for noise in Section 3. Finally, in Section 4, we apply the adjusted test statistics on our real data sets and compare the  $p$ -values obtained against those from the CSN test.

## 1. Reexamining Significance Testing for Empirical Distributions

Sometimes we have reasons to believe that our large data sets may be described by well known distributions, such as the normal distribution, power law distribution, exponential distribution, and so on, but with best-fit parameter values that we need to determine. The most common methods used to perform this *parameter estimation* are the Maximum Likelihood Estimation (MLE) (26), Maximum Entropy Method (MEM) (27–29), Regressions (30), and direct or indirect computation of moments (31). Since it is possible to fit any distribution to any data set, we need to compute its *goodness of fit*, which can be the KS distance (32), the coefficient of determination ( $R^2$ ) and other forms of distance measure (33, 34).

In a recent statement, the American Statistical Association warned the scientific community that the  $p$ -value “was never intended to be a substitute for scientific reasoning” (35, para. 2), and outline six principles that can prevent its misuse (36). A *Nature* commentary on this statement also added that “[r]esearchers should describe not only the data analyses that produced statistically significant results, . . . , but all statistical tests and choices made in calculations” (37, para. 3). We heed the warning in this paper, but argue that when properly computed and interpreted, the  $p$ -value is useful in that it provides a quantitative and objective alternative to visual inspection of the fits. The latter is frequently subjective and biased. This utility becomes important when we are comparing fits of two or more data sets to two or more distributions, and have the ambiguity of being able to choose from two or more definitions of goodness of fit. This is why we need to go beyond the goodness of fits, to establish how plausible different distributions are for different data sets.

In 2009, Clauset, Shalizi and Newman (CSN) did precisely this by coming up with a  $p$ -test model that use the well-known PL distribution as an illustration. They started by writing down the probability density function for PL distribution

$$f_{PL} = \frac{\alpha - 1}{x_{min}^{1-\alpha}} x^{-\alpha}, \quad [1]$$

for  $x \in [x_{min}, \infty)$ . The CSN  $p$ -test involves three major steps:

**CSN(i) MLE Estimation of  $\alpha$ .** Given an empirical data with  $S$  observations,  $\mathbf{Y} = \{y_1, y_2, \dots, y_S\}$  sorted such that  $y_i \leq y_{i+1}$ , the CSN algorithm (CSN(ia)) first constructs the  $S$  subsets  $\mathbf{X}^{(j)} = \{x_1^{(j)} = y_j, x_2^{(j)} = y_{j+1}, \dots, x_N^{(j)} = y_S\}$ . (CSN(ib)) For each  $\mathbf{X}^{(j)}$ , we estimate  $\alpha^{(j)}$  using the MLE method that maximizes the log-likelihood function,

$$\ln \mathbb{L}_{PL} = \ln \left[ \prod_{i=1}^N f_{PL}(x_i | \hat{\alpha}) \right] = N \ln \left( \frac{\hat{\alpha} - 1}{x_{min}} \right) - \hat{\alpha} \sum_{i=1}^N \frac{x_i}{x_{min}}. \quad [2]$$

Applying the maximizing condition  $\frac{\partial(\ln \mathbb{L})}{\partial \alpha} = 0$  yields,

$$\hat{\alpha} = 1 + \left\langle \ln \frac{x}{x_{min}} \right\rangle^{-1}, \quad [3]$$

where the hat indicates an estimated parameter.

**CSN(ii) KS Distance.** If  $\mathbf{X}$  follows probability distribution function  $f_X$  with cumulative distribution function  $F_X$ , then its probability integral transform  $u = F_X(x)$  is a standard uniform distribution function ( $U(0, 1)$ ). For any PL distributed sample  $\mathbf{X} = \{x_1 = x_{min}, x_2, \dots, x_N\}$  with estimated  $\hat{\alpha}$ , we (CSN(iia)) first transform the sample to  $U^{(s)} = \{u_i^{(s)} = F_{PL}(x_i | \hat{\alpha})\}_{i=1}^N$ . (CSN(iib)) Then we calculate the KS distance

$$d_{KS} = \forall_{i=1}^N \sup \left( \left| u_i^{(s)} - \frac{i}{N} \right| \right) \quad [4]$$

between  $U^{(s)}$  and  $U(0, 1)$ . Here we make use of the fact that the CDF of  $U(0, 1)$  is a linear function,  $F_U(u) = u$ .

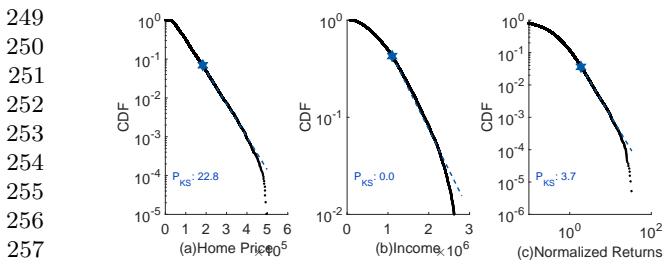
**CSN(iii) Determining  $x_{min}$ .** To determine  $x_{min}$ , (CSN(iii)) we calculate the KS distance for each  $\mathbf{X}^{(j)}$  with its corresponding  $\hat{\alpha}^{(j)}$ . (CSN(iiib)) The set  $\mathbf{X}^{(j)}$  that yields the lowest KS distance ( $d_{KS}^{(em)}$ ) gives us  $\hat{x}_{min}^{(em)} = y_j$  and  $\hat{\alpha}^{(em)} = \hat{\alpha}^{(j)}$ . The superscript ‘(em)’ indicates a parameter obtained from empirical data.

**CSN(iv) Significance Testing.** After  $\hat{\alpha}^{(em)}$  and  $\hat{x}_{min}^{(em)}$  have been estimated from  $\mathbf{Y} = \{y_1, y_2, \dots, y_S\}$ , we test how plausible it is for  $\mathbf{X} = \{x_1 = \hat{x}_{min}^{(em)}, x_2, \dots, x_N\} \subset \mathbf{Y}$  to be a sample taken from a PL distribution. This is done by (CSN(iv)) sampling the PL  $M$  times using  $\hat{\alpha}^{(em)}$  and  $\hat{x}_{min}^{(em)}$ . (CSN(ivb)) For the  $m^{th}$  simulated sample we go through CSN(i) to CSN(iii) to obtain  $d_{KS}^{(m)}$ . (CSN(ivc)) The significance measure

$$p = \frac{1}{M} \sum_{m=1}^M \mathbb{I}_{\{d_{KS}^{(em)} < d_{KS}^{(m)}\}}. \quad [5]$$

is the fraction of simulated samples whose fits are poorer than that of the data.

Extending the CSN method to other distributions, we performed  $p$ -testing on the Taiwan home price per square foot (fitted to EXP), Taiwan income (fitted to EXP), and the Straits Times Index normalized return (fitted to PL) (see SI Section 3 for more descriptions on the data sets). The fits and  $p$ -values are shown in Fig. 1. All fits are visually good yet only the  $p$ -value for the Taiwan housing is appreciable. We realized the reason for this is simple: while the EXP and PL distributions are defined over  $(0, \infty)$ , when we collect data from the real world we can only obtain a finite number of elements. Moreover, the largest element in the data is finite.



**Fig. 1.** *p*-testing on [left] 2012 to 2014 Taiwan home price per square foot in order of  $\times 10^5$  (fitted to EXP), [middle] 2012 Taiwan lower-tail income in million (fitted to EXP), and [right] 2009 to 2016 Strait Time Index normalized return (fitted to PL). The black dots represent empirical data while the magenta dashed line represents the fit. All fits are visually good, yet only the *p*-value ( $P_{KS}$  in percentage) for Taiwan home price is appreciable.

However, existing tests for statistical significance generally do not account for the effects produced by having a finite number of elements (FNE) and a finite largest element (FLE). In the next section we will explain how the parameters and test statistics can be adjusted for FNE and FLE.

## 2. Finite-Sample Adjustments

**Parameter Adjustment for Finite Largest Element.** Here, we will illustrate the effects of FLE using an asymptotic EXP distribution. The same discussion can be generalized to other distributions (see SI Section 1).

The EXP distribution is defined as

$$f_{EXP}(x) = \beta \exp[-\beta(x - x_{min})], \quad [6]$$

with  $\beta$  as a sole parameter for  $x \in [x_{min}, \infty)$ . Maximizing the likelihood function  $L = \prod_{i=1}^N P(X = x_i | x_{min}, \hat{\beta})$ , we find the estimated parameter

$$\hat{\beta} = \frac{1}{\langle x \rangle - x_{min}}. \quad [7]$$

If we use the mean obtained from data  $\langle x \rangle_{data}$  as  $\langle x \rangle$  in Eq. (7) we will obtain the unadjusted estimator  $\hat{\beta}_{unadj}$ . However due to the FLE, we can only average up till  $x_{max}$ . As such  $\langle x \rangle_{data}$  will be biased downwards and Eq. (7) over-estimates  $\hat{\beta}$ .

To adjust for the FLE, we add the truncated part back into  $\langle x \rangle_{data}$ , to define the adjusted  $\langle x \rangle_{adj}$  as

$$\begin{aligned} \langle x \rangle_{adj} &= \langle x \rangle_{data} \int_{x_{min}}^{x_{max}} f_{EXP}(x) dx + \int_{x_{max}}^{\infty} x f_{EXP}(x) dx. \\ &= \langle x \rangle_{data} \{1 - \exp[-\beta(x_{max} - x_{min})]\} + \\ &\quad \frac{\exp[-\beta(x_{max} - x_{min})]}{\beta} [\beta x_{max} + 1]. \end{aligned} \quad [8]$$

Inserting  $\langle x \rangle_{adj}$  into Eq. (7), we obtain a nonlinear equation

$$[\hat{\beta}_{adj} (\langle x \rangle_{data} - x_{min}) + 1] \exp[-\hat{\beta}_{adj} (x_{max} - x_{min})] + \hat{\beta}_{adj} (\langle x \rangle_{data} - x_{min}) - 1 = 0 \quad [9]$$

that we need to solve to obtain  $\hat{\beta}_{adj}$ . This is an extension of the series of work (38–40) on improving parameter estimation for PLs given both  $x_{min}$  and  $x_{max}$ .

To test the performance of this adjustment formula, we simulated 1000 sets of EXP distributed data for  $10^{-4} \leq \beta_T \leq 10^2$ , by using the inverse cumulative function for EXP distribution

$$F_{EXP}^{-1}(u, \beta_T) = x_{min} - \frac{1}{\beta_T} \ln(1 - u). \quad [10]$$

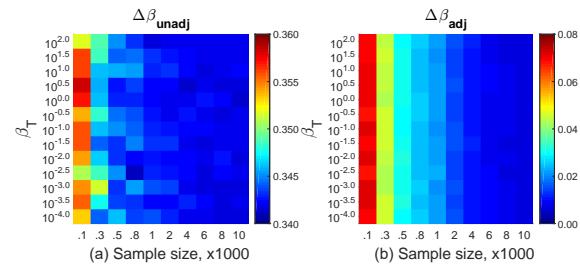
This transforms  $U(0, 1)$  distributed elements  $\{u_i\}$  to EXP distributed elements  $\{x_i\}$ . Using this transformation  $F_{EXP}^{-1}$ , 0 and 1 map to  $x_{min}$  and  $\infty$  respectively. It is also useful to note that equation Eq. (10) is the inverse of the CDF of the EXP distribution,

$$F_{EXP}(x, \beta_T) = 1 - \exp[\beta(x_{min} - x)]. \quad [11]$$

To simulate the effect of a FLE with  $x_{max} = F_{EXP}^{-1}(0.9)$ , we sampled 1000 sets of EXP distributed data using  $U(0, 0.9)$  instead of  $U(0, 1)$  with  $x_{min} = 0$ . Thereafter, we estimated  $\hat{\beta}_{unadj}$  and  $\hat{\beta}_{adj}$  using Eq. (7) and Eq. (9). Fig. 2 shows the relative estimation errors

$$\Delta\beta = \frac{\sqrt{\langle(\hat{\beta} - \beta_T)^2\rangle}}{\beta_T} \quad [12]$$

of  $\hat{\beta}_{unadj}$  and  $\hat{\beta}_{adj}$  with respect to the true beta  $\beta_T$ . As we can see from the Fig.,  $\Delta\hat{\beta}_{unadj}$  is about 38% for small samples  $N \sim 10^2$  and decreases to 34% for large samples  $N \sim 10^4$ . On the other hand,  $\Delta\hat{\beta}_{adj}$  starts at 20%, but decreases to 2% as the number of data points is increased. Although it can be shown that the  $\hat{\beta}_{unadj}$  is an unbiased estimator (24, 41), we find it converging very slowly with increasing sample size in the unfortunate situation of a small  $x_{max}$ . In contrast,  $\hat{\beta}_{adj}$  converges very quickly even for small  $x_{max}$  as we have accounted for the FLE.



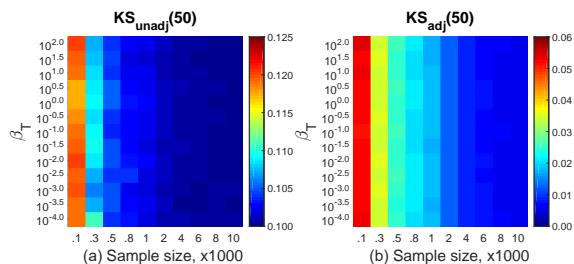
**Fig. 2.** Relative estimation errors of (a)  $\hat{\beta}_{unadj}$  and (b)  $\hat{\beta}_{adj}$  measured from 1000 simulated samples using different  $\beta_T$  and  $N$  with  $x_{min} = 0$  and  $x_{max} = F_{EXP}^{-1}(0.9)$ . Due to the FLE,  $\Delta\hat{\beta}_{unadj}$  remains high (close to the theoretical relative error of  $\epsilon(\delta = 0.1, x_{min} = 0) = 0.1[1 - \ln(0.1)] \approx 33\%$ ) even for large  $N$ . In contrast,  $\Delta\hat{\beta}_{adj}$  decreases rapidly with increasing  $N$ .

In the SI Section 1, we show details for our derivation of the theoretical relative error

$$\begin{aligned} \hat{\beta} &\approx \beta_T + \beta_T [\beta_T x_{max} + 1] \exp(-\beta_T(x_{max} - x_{min})) + \\ &\quad \mathcal{O}\{(\beta_T x_{max} + 1)^2 \exp(-2\beta_T(x_{max} - x_{min}))\}. \end{aligned} \quad [13]$$

Eq. (13) shows the independence of the estimation error on sample size. This tells us that the  $\hat{\beta}_{unadj}$  is always larger than the  $\beta_T$  because of the FLE effect. The convergence rate then depends on how rapidly  $x_{max}$  approaches infinity with increasing sample size.

373 **Test Statistic Adjustment for FLE.** For a finite sample,  
 374  $F_{EXP}(x) < 1$  for all  $x < \infty$ . Mathematically, this means  
 375 that  $F_{EXP}(x) \sim U(0, 1 - \delta)$ , where  $F_{EXP}^{-1}(x_{max}) = 1 - \delta$ .  
 376 This observation is important, because  $d_{KS}$  is obtained by  
 377 comparing  $U^{(s)} = \{u_i^{(s)} = F_{EXP}(x_i|\hat{\beta})\}_{i=1}^N$  against  $U(0, 1)$   
 378 (see Eq. (4)). This tell us that for a fair comparison, we need  
 379 to rescale all elements in  $U^{(s)}$  by a factor of  $1/(1 - \delta)$ . Fig. 3  
 380 shows the  $d_{KS}$  measured for the 1000 sets of EXP distributed  
 381 data with finite largest element  $x_{max} = F^{-1}(0.9)$  for various  
 382  $\beta_T$  and sample sizes  $N$ . For each sample, we use Eq. (9) to  
 383 estimate the  $\hat{\beta}_{adj}$  and transformed this data to  $U^{(s)}$  using  
 384 Eq. (11). After that, we measure  $d_{KS}$  with Eq. (4) to obtain  
 385 unadjusted KS distance,  $KS_{unadj}$  and adjusted KS distance,  
 386  $KS_{adj}$  using the non-rescaled and rescaled  $U^{(s)}$  respectively.  
 387  $KS_{unadj}$  goes from 0.14 for small samples  $N \sim 10^2$ , to 0.10  
 388 for large samples  $N \sim 10^5$ . In contrast,  $KS_{adj}$  decrease from  
 389 0.06 for small samples to 0.006 for large samples. Last but not  
 390 least, it is good to note that the  $d_{KS}$  measured is independent  
 391 of  $\beta_T$ .

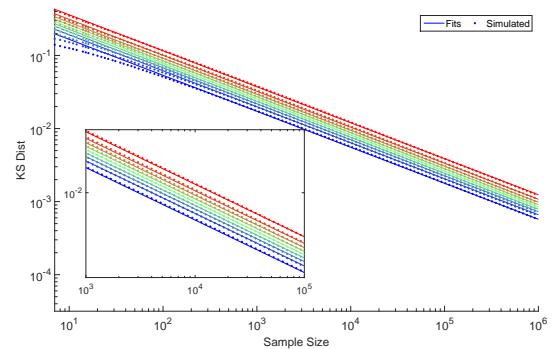


403 **Fig. 3.** The median KS distances for (a)  $KS_{unadj}$  and (b)  $KS_{adj}$  measured from  
 404 1000 simulated samples using different  $\beta_T$  and  $N$ . The  $x_{min}$  is set to 0 and  
 405  $x_{max} = F_{EXP}^{-1}(0.9)$ . Because of the FLE,  $KS_{unadj}$  remains above  $\delta = 0.10$   
 406 while  $KS_{adj}$  converges to zero for large  $N$ .

410 **Adjustment for Finite Number of Elements.** Until now, we have  
 411 only discussed adjustments to the estimated parameter and the  
 412 KS distance to eliminate the bias caused by the FLE. Besides  
 413 the FLE effect, we also need to consider the bias caused by  
 414 having a finite number of elements in the sample. As we can  
 415 see from Fig.s 3 , the KS distance decreases as the sample size  
 416 increases. Therefore, in order to have a fair comparison of the  
 417 goodness of fit for various sample sizes, we need to determine  
 418 how  $d_{KS}$  changes as a function of  $N$ . To do this, we simulated  
 419  $10^6$  samples of various sizes  $N$  from  $U(0, 1)$ . For each sample  
 420 we determined  $d_{KS}$  using Eq. (4), so that for each  $N$  we end  
 421 up with  $10^6$  KS distances. In Fig. 4 we show the KS distances  
 422 at different deciles, which exhibits the asymptotic behavior

$$d_{KS}(\varphi_{KS}, N) = \frac{\left(\frac{100}{\varphi_{KS}} - 1\right)^{-0.176} \exp(-0.274)}{N^{0.492}}, \quad N > 50 \quad [14]$$

423 that we settled for, after experimenting with several functional  
 424 forms. This result agrees with our expectation that  $d_{KS} \rightarrow 0$   
 425 as  $N \rightarrow \infty$ . It also suggests that if we have two samples  
 426 with sizes  $N_1$  and  $N_2$  from the same distribution, we should  
 427 compare  $N_1^{0.492} d_{KS}^{(1)}$  against  $N_2^{0.492} d_{KS}^{(2)}$ . Otherwise, if  $N_2 > N_1$   
 428 then naturally  $d_{KS}^{(2)} < d_{KS}^{(1)}$  and we will be lead to the wrong  
 429 conclusion that the  $N_2$  sample fits the distribution better.



449 **Fig. 4.** Log-log plot of  $d_{KS}$  against  $N$  for different deciles going from the 10<sup>th</sup>  
 450 percentile (blue) to the 90<sup>th</sup> (red), obtained from  $10^6$  simulations.

### 3. The Effects of Random Noise

451 Besides having to work with finite samples and finite largest  
 452 elements, we will also in practice encounter imperfections while  
 453 collecting samples for various reasons, such as undetected sam-  
 454 ples, contamination by background noise, and recording errors.  
 455 We call such noises that occur at the element level *elementary*  
 456 noise. When we convert these samples to a distribution, noise  
 457 will also be present at the distribution level that we refer to as  
 458 *distribution noise*. In principle the information at the distribu-  
 459 tion level is more robust compared to the elementary level, as  
 460 we expect random and thus uncorrelated noise to cancel each  
 461 other. This means that the distribution is less sensitive to  
 462 elementary noise, but we still worry whether the distribution  
 463 noise may play an important role in our test of statistical  
 464 significance. In order to account for the effects of distribution  
 465 noise, we need to first be able to quantify distribution noise,  
 466 and thereafter understand how it affects significance testing.

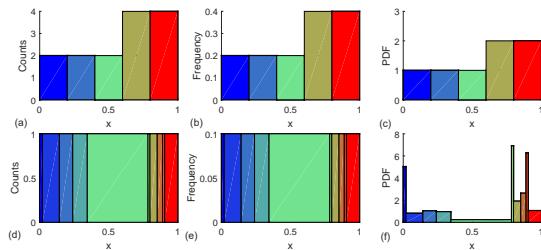
467 Suppose we now randomly generate a set of EXP data.  
 468 After adjusting for FLE, we obtained the distribution  
 469 parameters and use it to transform this set to  $U^{(s)} = \{u_i^{(s)} =$   
 470  $F_{EXP}(x_i|\hat{\beta})\}_{i=1}^N$  following the procedure outlined in Section  
 471 2. Then as illustrated in Fig. 5(a)-(b), a natural way to mea-  
 472 sure the distribution noise is to plot the histogram, count  
 473 the frequency for each bin, and compare it to the expected  
 474 frequency from  $U(0, 1)$ . Since this can be more accurately  
 475 done for smaller bin sizes, we use the intervals between sorted  
 476 elements as a collection of non-uniform bins, as shown in  
 477 Fig. 5(c)-(d). For a data set consisting of  $N$  elements, each bin  
 478 carry a weight of  $1/N$ , evenly distributed within the interval  
 479  $(u_{i-1}, u_i)$ , such that the probability density is

$$f(u_{i-1}, u_i) = \frac{\frac{1}{N}}{u_i - u_{i-1}}. \quad [15]$$

480 As the theoretical probability density for  $U(0, 1)$  is 1, we define  
 481 the distribution noise  $d_{DN}$  mathematically to be

$$\begin{aligned} d_{DN} &= \sqrt{\frac{\sum_{i=1}^N (u_i - u_{i-1})^2 [f(u_{i-1}, u_i) - 1]^2}{\sum_{i=1}^N (u_i - u_{i-1})^2}} \\ &= \sqrt{\frac{\sum_{i=1}^N (u_i - u_{i-1})^2 \left(\frac{1}{N(u_i - u_{i-1})} - 1\right)^2}{\sum_{i=1}^N (u_i - u_{i-1})^2}}, \end{aligned} \quad [16]$$

497 where  $u_0 = 0$  and  $u_N = 1$ . We need to weigh the deviation  
 498 of each bin by  $(u_i - u_{i-1})^2$  because the bins are non-uniform.  
 499 and to keep  $d_{DN}$  finite.



500  
 501 **Fig. 5.** Illustration of the distribution noise we would measure if we sample 10  
 502 elements from  $U(0, 1)$ , rescaled such that the largest element becomes 1. In (a)-  
 503 (b) we use 5 uniform bins whereas in (c)-(d) we use the intervals between sorted  
 504 elements as the bins. The counts are shown on the left, whereas the probability  
 505 densities calculated using Eq. (15) are shown on the right.  
 506  
 507

508 **Relation between Distribution Noise and Sample Size.** As  
 509 with Section 2, we simulated  $10^6$  samples from  $U(0, 1)$  with  
 510 different  $N$ . For each sample, we calculate the distribution  
 511 noise  $d_{DN}$  using Eq. (16) and plot its deciles against  $N$  as  
 512 shown in Fig. 6. After experimenting with several functional  
 513 forms, we write down the relationship between  $d_{DN}$  and  $N$  at  
 514 percentile  $\varphi_{DN}$  as

$$515 \quad d_{DN}(\varphi_{DN}, N) = d_{DN}^{(*)} + \\ 516 \quad \Phi(\varphi_{DN} - 50) \frac{\exp\left(-\frac{[50 - |\varphi_{DN} - 50|]^{0.430}}{|\varphi_{DN} - 50|^{0.302}}\right)}{N^{0.495}}, \quad [17]$$

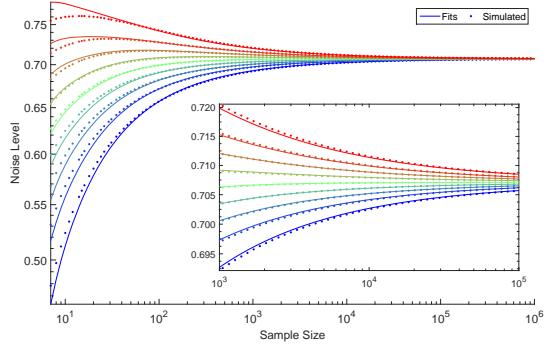
517 where  $\Phi(x)$  represents the sign of  $x$ , and

$$518 \quad d_{DN}^{(*)} = \sqrt{\frac{1}{2} + \frac{2-N}{2N^2}} \left( \frac{N}{N+0.5} \right) \quad [18]$$

519 is the analytically derived distribution noise (See SI Section  
 520 2), that converges to  $1/\sqrt{2}$  as  $N \rightarrow \infty$ . This result suggests  
 521 that if we have two samples with sizes  $N_1$  and  $N_2$  with  $N_2 > N_1$   
 522 from the same distribution, we should compare  $N_1^{0.495}(d_{DN}^{(1)} -$   
 523  $1/\sqrt{2})$  against  $N_2^{0.495}(d_{DN}^{(2)} - 1/\sqrt{2})$ . Otherwise, we risk making  
 524 wrong conclusion that the  $N_2$  sample fits the distribution  
 525 better if  $d_{DN}^{(1)} > d_{DN}^{(2)}$ .

526 **Relationship between Distribution Noise and KS Distance.** As  
 527 measures for statistical deviations,  $d_{DN}$  and  $d_{KS}$  are different  
 528 in that  $d_{DN}$  measures deviation at the probability density  
 529 level, whereas the  $d_{KS}$  measure it at the cumulative density  
 530 level. As a result,  $d_{KS}$  assigns more weight to the tail of the  
 531 distribution, while  $d_{DN}$  is more sensitive to deviations in the  
 532 body of the distribution. Therefore, if we wish to combine  
 533 these two measures to estimate the significance level, we need  
 534 to first investigate the relationship between  $d_{KS}$  and  $d_{DN}$ . We  
 535 do this by simulating  $10^6$  samples from  $U(0, 1)$  for various  
 536 sample sizes, and for each sample, we calculate  $d_{KS}$  and  $d_{DN}$   
 537 using Eq. (4) and Eq. (16) respectively, to obtain  $10^6$  pairs  
 538 of  $d_{KS}$  and  $d_{DN}$ . We then compute the Pearson correlation  
 539 between  $d_{KS}$  and  $d_{DN}$  and learned that (see SI Section 2  
 540 for the comparison of fits)

$$541 \quad \rho_{d_{KS}, d_{DN}}(N) = \frac{e}{N^{0.481}}. \quad [19]$$



541 **Fig. 6.** Relationship between distribution noise  $d_{DN}$  and sample size  $N$  at deciles  
 542 going from the 10<sup>th</sup> percentile (blue) to the 90<sup>th</sup> (red), obtained from  $10^6$  simulations.  
 543 The  $d_{DN}$  value converges to  $1/\sqrt{2}$  as  $N$  increases.  
 544  
 545

546 As expected,  $d_{KS}$  is positively correlated with  $d_{DN}$ . Since  $d_{KS}$   
 547 is a measure at the cumulative level, the random distribution  
 548 noises cancel each other, thus the correlation between  $d_{KS}$   
 549 and  $d_{DN}$  vanishes as  $N \rightarrow \infty$ .

#### 4. Application to Significance Testing

550 **Significance Level for a Given Distribution.** To perform significance  
 551 testing given  $d_{KS}$  and  $d_{DN}$ , we need the percentile  
 552 values  $\varphi_{KS}$  and  $\varphi_{DN}$ .  $\varphi_{KS}$  can be obtained by inverting  
 553 Eq. (14), as

$$554 \quad \varphi_{KS}(d_{KS}, N) = \frac{100}{\left(1 + (d_{KS} N^{0.492} \exp(0.274))^{-\frac{1}{0.176}}\right)}. \quad [20]$$

555 Similarly, we invert Eq. (17), and solve

$$556 \quad \varphi_{DN}^{0.430} + (50 - \varphi_{DN})^{0.302} \ln(|\eta| N^{0.495}) = 0, \quad \eta < 0 \\ 557 \quad \varphi_{DN} - 50 = 0, \quad \eta = 0$$

$$558 \quad (100 - \varphi_{DN})^{0.430} + (\varphi_{DN} - 50)^{0.302} \ln(|\eta| N^{0.495}) = 0, \quad \eta > 0 \quad [21]$$

559 to get  $\varphi_{DN}$ , where  $\eta = d_{DN} - \langle d_{DN} \rangle$ .

560 Substituting the empirical KS distance  $d_{KS}^{(em)}$  and empirical  
 561 distribution noise  $d_{DN}^{(em)}$  into Eq. (20) and Eq. (21), we obtain  
 562  $\varphi_{KS}^{(em)}$  and  $\varphi_{DN}^{(em)}$ . This is an alternative way of obtaining the  
 563 p-value without the need to perform Monte-Carlo (re)sampling  
 564 again (CSN method), since we have already done so in Section  
 565 2 and 3. The percentage of simulated  $U(0, 1)$  samples with  
 566  $d_{KS/DN} > d_{KS/DN}^{(em)}$  is  $100 - \varphi_{KS/DN}^{(em)}$ . Since  $d_{KS}$  and  $d_{DN}$   
 567 are not independent (Eq. (19)), we discount the correlation  
 568 between  $d_{KS}$  and  $d_{DN}$ , and define the significance level ( $p$ -  
 569 value) as

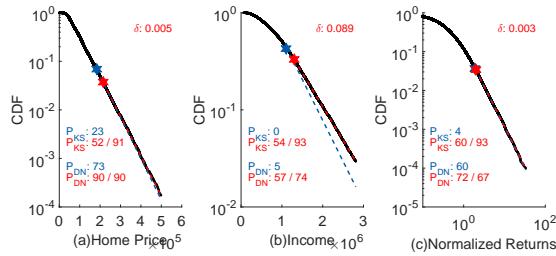
$$570 \quad p(\varphi_{KS}, \varphi_{DN}, N) = \sqrt{\left(1 - \frac{\varphi_{KS}}{100}\right) \left(1 - \frac{\varphi_{DN}}{100}\right) \left(1 - \frac{e}{N^{0.481}}\right)}, \quad [22]$$

571 to avoid overestimating the significance level.

572 **Fitting to Empirical Data.** We follow the steps outlined in the  
 573 CSN algorithm (Section 1) to fit the empirical data, but with  
 574 two important modifications: (Ii) the parameters (**CSN(ib)**)  
 575 and goodness of fit (**CSN(ibb)**) are adjusted for the finite  
 576 largest element; and (Iii) the p-value (**CSN(ivc)**) is adjusted

for the finite number of elements effect. Meanwhile, optional modifications are (O<sub>i</sub>) to incorporate distribution noise as another dimension for goodness of fit, so that the *p*-value can be determined via  $d_{KS}^{(em)}$ ,  $d_{DN}^{(em)}$ , or both; (O<sub>ii</sub>) instead of using bootstrapping to determine the *p*-value in the CSN method, which is very slow for large samples, one can use the fast inversion formulae Eq. (20), Eq. (21), or Eq. (22).

Fig. 7 shows the fits and *p*-testing results for Taiwan housing price, Taiwan wealth, and Straits Times Index normalized returns. It reassures that after modifications the *p*-values of all distributions increased. In particular, the two distributions (Fig. 7(b),(c)) that did not meet the  $p > 0.1$  criterion (as suggested by (24)) before modification, now have  $p > 0.5$ . This is in agreement with our visual assessment of the three fits. We also understand now that a large  $\delta$  (small  $x_{max}$ ) is the main reason for Taiwan wealth to fail *p*-testing before adjustment (although the fit is visually good). In general, our correction formulas perform the best when  $\delta$  is large due to small sample sizes or truncations. Readers can refer to SI Section 4 for more plots and instances where small  $\delta$  values affects the significance testing.



**Fig. 7.** *p*-testing results for (a) 2012 to 2014 Taiwan home price per square foot in order of  $\times 10^5$  (fitted to EXP), (b) 2012 Taiwan lower-tail income in millions (fitted to EXP), and (c) 2009 to 2016 Straits Times Index normalized return (fitted to PL) before and after finite-sample adjustments. In the Fig., the empirical CDF that is adjusted for FLE is shown as black dots, while the unadjusted and adjusted fits are shown as magenta and red dashed line respectively.  $P_{KS/DN}$ -values (in percentage) are for unadjusted (magenta) and adjusted (red) fits. We separate the *p*-values obtained using the CSN method (left) from those using Eq. (20) or Eq. (21) (right) by a '/'.

There are several limitations one should note while obtaining  $P_{KS/DN}$  using Eq. (20) or Eq. (21). First, it is only applicable to large samples (see Figs. 4 and 6). Second, these equations are obtained after experimenting with several functional forms and are only approximate. Lastly,  $p_{KS}$  measured using the CSN method are consistently smaller than based on Eq. (20). This is due to the CSN algorithm having an extra step to select  $x_{min}$  that minimizes  $d_{KS}$  of each simulated sample, and thus the algorithm is stricter than our fast inversion formulae. However the inversion formulae Eq. (20) and Eq. (21) is convenient and provides an upper bound for  $P_{KS/DN}$ . We make the codes for these two methods available at <https://github.com/darreltay/p-Testing> but leave it to the reader to decide which method to use.

All in all, when we test for statistical significance, we need to be aware of finite sample effects, namely the finite largest element effect and the finite number of elements effect. Beyond the KS distance measured at the cumulative distribution level, we also introduce an alternative measure of the goodness of fit based on the distribution noise at the probability density level.

**ACKNOWLEDGMENTS.** The authors would like to thank Chou Chung-I for directing us to the Taiwanese data sets.

1. Kac M, Kiefer J, Wolfowitz J (1955) On tests of normality and other tests of goodness of fit based on distance methods. *Ann. Math. Statist.* 26(2):189–211.
2. D'Agostino RB (1970) Transformation to normality of the null distribution of g1. *Biometrika* pp. 679–681.
3. Jarque CM, Bera AK (1987) A test for normality of observations and regression residuals. *International Statistical Review/Revue Internationale de Statistique* pp. 163–172.
4. Shapiro S, Wilk M (1965) An analysis of variance test for normality. *Biometrika* 52(3):591–611.
5. Anderson TW, Darling DA (1952) Asymptotic theory of certain "goodness of fit" criteria based on stochastic processes. *Ann. Math. Statist.* 23(2):193–212.
6. Anderson TW, Darling DA (1954) A test of goodness of fit. *Journal of the American Statistical Association* 49(268):765–769.
7. Jr. FJM (1951) The kolmogorov-smirnov test for goodness of fit. *Journal of the American Statistical Association* 46(253):68–78.
8. Lilliefors HW (1967) On the kolmogorov-smirnov test for normality with mean and variance unknown. *Journal of the American Statistical Association* 62(318):399–402.
9. Razali NM, Wah YB, et al. (2011) Power comparisons of shapiro-wilk, kolmogorov-smirnov, lilliefors and anderson-darling tests. *Journal of statistical modeling and analytics* 2(1):21–33.
10. Newman ME (2005) Power laws, pareto distributions and zipf's law. *Contemporary physics* 46(5):323–351.
11. Mantegna RN, Stanley HE (1995) Scaling behaviour in the dynamics of an economic index. *Nature* 376(6535):46–49.
12. Pierrou V, Gopikrishnan P, Amaral LAN, Meyer M, Stanley HE (1999) Scaling of the distribution of price fluctuations of individual companies. *Physical review e* 60(6):6519.
13. Gopikrishnan P, Pierrou V, Amaral LAN, Meyer M, Stanley HE (1999) Scaling of the distribution of fluctuations of financial market indices. *Physical Review E* 60(5):5305.
14. Teh BK, Cheong SA (2016) The asian correction can be quantitatively forecasted using a statistical model of fusion-fission processes. *PLoS one* 11(10):e0163842.
15. Zipf GK (1949) Human behavior and the principle of least effort. *Ed: Addison-Wesley, Reading, MA.*
16. Cancho RF, Solé RV (2001) The small world of human language. *Proceedings of the Royal Society of London B: Biological Sciences* 268(1482):2261–2265.
17. Auerbach F (1913) Das gesetz der bevölkerungskonzentration. *Petermanns Geographische Mitteilungen* 59:74–76.
18. Gabai X, Ioannidis YM (2004) The evolution of city size distributions. *Handbook of regional and urban economics* 4:2341–2378.
19. MacKay N (2010) London house prices are power-law distributed. *arXiv preprint arXiv:1012.3039*.
20. Ohnishi T, Mizuno T, Shimizu C, Watanabe T (2012) Power laws in real estate prices during bubble periods. *International Journal of Modern Physics: Conference Series* 16:61–81.
21. Tay DJ, Chou CI, Li SP, Tee SY, Cheong SA (2016) Bubbles are departures from equilibrium housing markets: Evidence from singapore and taiwan. *PLoS one* 11(11):e0166004.
22. Mandelbrot B (1960) The pareto-levy law and the distribution of income. *International Economic Review* 1(2):79–106.
23. Yakovenko VM, Rosser JR (2009) Colloquium: Statistical mechanics of money, wealth, and income. *Reviews of Modern Physics* 81(4):1703.
24. Clauset A, Shalizi CR, Newman ME (2009) Power-law distributions in empirical data. *SIAM review* 51(4):661–703.
25. Brzezinski M (2014) Do wealth distributions follow power laws? evidence from 'rich lists'. *Physica A: Statistical Mechanics and its Applications* 406:155–162.
26. Fisher RA (1912) On an absolute criterion for fitting frequency curves. *Messenger of Mathematics* 41:155–160.
27. Kumphon B (2012) Maximum entropy and maximum likelihood estimation for the three-parameter kappa distribution. *Open Journal of Statistics* 2(4):415–419.
28. Hradil Z, Reháček J (2006) Likelihood and entropy for statistical inversion in *Journal of Physics: Conference Series*. (IOP Publishing), Vol. 36, p. 55.
29. Akaike H (1998) *Information Theory and an Extension of the Maximum Likelihood Principle*, eds. Parzen E, Tanabe K, Kitagawa G. (Springer New York, New York, NY), pp. 199–213.
30. Bates DM, Watts DG (1988) *Nonlinear regression analysis and its applications*. (Wiley).
31. Wooldridge JM (2001) Applications of generalized method of moments estimation. *The Journal of Economic Perspectives* 15(4):87–100.
32. Jr. FJM (1951) The kolmogorov-smirnov test for goodness of fit. *Journal of the American Statistical Association* 46(253):68–78.
33. Cameron AC, Windmeijer FA (1997) An r-squared measure of goodness of fit for some common nonlinear regression models. *Journal of Econometrics* 77(2):329 – 342.
34. Janczura J, Weron R (2012) Black swans or dragon-kings? a simple test for deviations from the power law. *The European Physical Journal Special Topics* 205(1):79–93.
35. American Statistical Association (2016) ASA *p*-value statement viewed > 150,000 times (<https://www.amstat.org/ASA/News/ASA-P-Value-Statement-Viewed-150000-Times.aspx>). Accessed: 2017-03-07.
36. Wasserstein RL, Lazar NA (2016) The asa's statement on *p*-values: context, process, and purpose. *Am Stat* 70(2):129–133.
37. Baker M (2016) Statisticians issue warning over misuse of *p* values. *Nature* 531:151.
38. Alstott J, Bullmore E, Plenz D (2014) powerlaw: a python package for analysis of heavy-tailed distributions. *PLoS one* 9(1):e85777.
39. Yu S, Klaus A, Yang H, Plenz D (2014) Scale-invariant neuronal avalanche dynamics and the cut-off in size distributions. *PLoS one* 9(6):e99761.
40. Marshall N, et al. (2016) Analysis of power laws, shape collapses, and neural complexity: new techniques and matlab support via the ncc toolbox. *Frontiers in Physiology* 7.
41. Pitman EJ, Pitman EJJ (1979) *Some basic theory for statistical inference*. (Chapman and Hall London) Vol. 7.

745	42. Pyke R (1965) Spacings. <i>Journal of the Royal Statistical Society. Series B (Methodological)</i>	pp. 395–449.	807
746			808
747			809
748			810
749			811
750			812
751			813
752			814
753			815
754			816
755			817
756			818
757			819
758			820
759			821
760			822
761			823
762			824
763			825
764			826
765			827
766			828
767			829
768			830
769			831
770			832
771			833
772			834
773			835
774			836
775			837
776			838
777			839
778			840
779			841
780			842
781			843
782			844
783			845
784			846
785			847
786			848
787			849
788			850
789			851
790			852
791			853
792			854
793			855
794			856
795			857
796			858
797			859
798			860
799			861
800			862
801			863
802			864
803			865
804			866
805			867
806			868