

# How I Learned to Stop Worrying and Love The Central Limit Theorem

ChangHoon Hahn, Florian Beutler, Manodeep Sinha, Andreas Berlind  
*Lawrence Berkeley National Laboratory, 1 Cyclotron Rd, Berkeley CA 94720, USA*  
changhoon.hahn@lbl.gov

DRAFT --- e53b0f3 --- 2017-11-27 --- NOT READY FOR DISTRIBUTION

## ABSTRACT

abstract here

*Subject headings:* methods: statistical — galaxies: statistics — methods: data analysis — cosmological parameters — cosmology: observations — large-scale structure of universe

## 1. Introduction

- Talk about the use of Bayesian parameter inference and getting the posterior in LSS cosmology
- Explain the two major assumptions that go into evaluating the likelihood
- Emphasize that we are not talking about non-Gaussian contributions to the likelihood
- Emphasize the scope of this paper is to address whether one of the assumptions matters for galaxy clustering analyses.

## 2. Gaussian Likelihood Assumption

- Depending on Hogg’s paper maybe a simple illustration of how the likelihood assumption

## 3. Mock Catalogs

Mock catalogs are play an indispensable role in standard cosmological analyses of LSS studies. They’re used for testing analysis pipelines (Beutler et al. 2017; Grieb et al. 2017; Tinker & et al. in preparation), testing the effect of systematics (Guo et al. 2012; Vargas-Magaña et al. 2014; Hahn et al. 2017; Pinol et al. 2017; Ross et al. 2017), and, most relevantly for this paper, estimating the covariance matrix (Parkinson et al. 2012; Kazin et al. 2014; Grieb et al. 2017; Alam et al. 2017; Beutler et al. 2017; Sinha et al. 2017). In fact, nearly

all current state-of-the-art LSS analyses use covariance matrices estimated from mocks to evaluate the likelihood for parameter inference.

While some argue for analytic estimates of the covariance matrix (e.g. [Mohammed et al. 2017](#)) or estimates directly from data by subsampling (e.g. [Norberg et al. 2009](#)), covariance matrices estimated from mocks have a number of advantages. Mock catalogs allow us to incorporate detailed systematic errors present in the data and variance beyond the volume of the data. Even for analytic estimates, large ensembles of mocks are crucial for validation (Slepian et al. 2017). Moreover, as we show later in this paper, mock catalogs allow us to quantify the non-Gaussianity of the likelihood and construct a more accurate estimate of the true likelihood.

In this paper, we focus on two LSS analyses: the powerspectrum multipole full shape analysis of [Beutler et al. \(2017\)](#) and group multiplicity function analysis of [Sinha et al. \(2017\)](#). Throughout the paper we will make extensive use of same the mock catalogs used in these analyses. Below, in this section, we give a brief description of these mocks.

### 3.1. MultiDark-PATCHY Mock Catalog

In their powerspectrum multipole full shape analysis, [Beutler et al. \(2017\)](#) use the MultiDark-PATCHY mock catalogs from [Kitaura et al. \(2016\)](#). These mocks are generated using the PATCHY code ([Kitaura et al. 2014, 2015](#)). They rely on large-scale density fields generated using augmented Lagrangian Perturbation Theory (ALPT [Kitaura & Heß 2013](#)) on a mesh. This mesh is then populated with galaxies based on a combined non-linear deterministic and stochastic biases. The mocks from the PATCHYcode are then calibrated to reproduce the galaxy clustering in the high-fidelity BigMultiDark  $N$ -body simulation ([Rodríguez-Torres et al. 2016; Klypin et al. 2016](#)).

The galaxies are then assigned stellar masses using the HADRON code ([Zhao et al. 2015](#)). And the SUGAR code ([Rodríguez-Torres et al. 2016](#)) is applied to combine different boxes, incorporate selection effects and masking to produce mock light-cone galaxy catalogs. The statistics of the resulting mocks are then compared to observations and the process is iterated to reach desired accuracy. We refer readers to [Kitaura et al. \(2016\)](#) for further details.

In total, [Kitaura et al. \(2016\)](#) generated a 12,228 mock light-cone galaxy catalogs for BOSS Data Release 12: 2048 for each southern and northern galactic caps of LOWZ, CMASS, combined samples. In [Beutler et al. \(2017\)](#), they use 2045 and 2048 for the northern galactic cap (NGC) and southern galactic cap (SGC) of the LOWZ+CMASS combined sample. [Beutler et al. \(2017\)](#) excluded 3 mock realizations, due to notable issues, which have been since been addressed. Therefore, in our analysis we use all 2048 mocks for both the NGC and SGC of the LOWZ+CMASS combined sample.

### 3.2. [Sinha et al. \(2017\)](#) Mocks

#### 3.3. $\mathbf{X}^{\text{mock}}$

$$\begin{aligned}\mathbf{X}_i &= [P_0(k)_i, P_2(k)_i, P_4(k)_i] \\ \mathbf{X}_i &= \zeta(N)_i \\ \mathbf{X} &= \{\mathbf{X}_i\}\end{aligned}$$

## 4. Quantifying the Likelihood non-Gaussianity

In [Sellentin et al. \(2017\)](#) **discuss how sellentin and hartlap’s stuff is an attempt to quantify the divergence**

A more direct approach can be taken to quantify the non-Gaussianity of the likelihood. We can calculate the divergence between the distribution of our observable,  $p(x)$ , and  $q(x)$  a multivariate Gaussian described by the average of the mocks and the covariance matrix  $\mathbf{C}$ . The following are two of the most commonly used divergences: the Kullback-Leibler (hereafter KL) divergence

$$D_{KL}(p \parallel q) = \int p(x) \log \frac{p(x)}{q(x)} dx \quad (1)$$

and the Rényi- $\alpha$  divergence

$$D_{R-\alpha}(p \parallel q) = \frac{1}{\alpha - 1} \log \int p^\alpha(x) q^{1-\alpha}(x) dx. \quad (2)$$

In the limit as  $\alpha$  approaches 1, the Rényi- $\alpha$  divergence is equivalent to the KL divergence.

Of course, in our case, we don’t know  $p(x)$  — *i.e.* the distribution of our observable. If we did, we would simply use that instead of bothering with the covariance matrix or this paper. We can, however, still estimate the divergence using nonparametric estimators ([Wang et al. 2009](#); [Póczos et al. 2012](#); [Krishnamurthy et al. 2014](#)). These estimators, allows us to estimate the divergence directly from samples  $X_{1:n} = \{X_1, \dots, X_n\}$  and  $Y_{1:m} = \{Y_1, \dots, Y_m\}$  drawn from  $p$  and  $q$  respectively:  $\hat{D}_\alpha(X_{1:n} \parallel Y_{1:m})$ .

For instance, the estimator presented in [Póczos et al. \(2012\)](#) allows us to estimate the kernel function of the Rényi- $\alpha$  divergence,

$$D_\alpha(p \parallel q) = \int p^\alpha(x) q^{1-\alpha}(x) dx. \quad (3)$$

using  $k$ th nearest neighbor density estimators. Let  $\rho_k(x)$  denote the Euclidean distance of the  $k$ th nearest neighbor of  $x$  in the sample  $X_{1:n}$  and  $\nu_k(x)$  denote the Euclidean distance of the  $k$ th nearest neighbor of  $x$  in the sample  $Y_{1:m}$ . Then  $D_\alpha(p \parallel q)$  can be estimated as

$$\hat{D}_\alpha(X_{1:n} \parallel Y_{1:m}) = \frac{B_{k,\alpha}}{n} \left( \frac{n-1}{m} \right)^{1-\alpha} \sum_{i=1}^n \left( \frac{\rho_k^d(X_i)}{\nu_k^d(X_i)} \right)^{1-\alpha}, \quad (4)$$

where  $B_{k,\alpha} = \frac{\Gamma(k)^2}{\Gamma(k-\alpha+1)\Gamma(k+\alpha-1)}$ . [Póczos et al. \(2012\)](#) goes to further prove that this estimated kernel function is asymptotically unbiased,

$$\lim_{n,m \rightarrow \infty} \mathbb{E}[\hat{D}_\alpha(X_{1:n} \parallel Y_{1:m})] = D_\alpha(p \parallel q). \quad (5)$$

Plugging  $\hat{D}_\alpha(X_{1:n} \parallel Y_{1:m})$  into Eq. 2, we get an estimator for the Rényi- $\alpha$  divergence. A similar estimator ([Wang et al. 2009](#)) can also be derived for the KL divergence (Eq. 1). We note that while the divergence estimators converge to the true divergence with a large enough samples, with a limited number of samples from the distribution, the estimators are noisy.

These divergence estimates have been applied to Support Vector Machines and used extensively in the machine learning and astronomical literature with great success **elaborate a lot more**

- Compile papers that use this divergence, [Ntampaka et al. \(2015, 2016\)](#)

For more details on the non-parametric divergence estimators, we refer readers to [Póczos et al. \(2012\)](#) and [Krishnamurthy et al. \(2014\)](#).

Now we can use the divergence estimators above to quantify the non-Gaussianity of the likelihood. More specifically, we're intersted in the divergence between the distribution  $p(x)$  sampled by the mock observables ( $\mathbf{X}^{\text{mock}}$ ) and the multivariate Gaussian distribution assumed in standard analyses, which is solely described by the mean and covariance calculated from the mocks ( $\mathcal{N}(\overline{\mathbf{X}^{\text{mock}}}, \mathbf{C})$ ). Since  $\mathbf{X}^{\text{mock}}$  is a sample from  $p(x)$ , we draw a reference sample  $\mathbf{Y}^{\text{ref}}$  from  $\mathcal{N}(\overline{\mathbf{X}^{\text{mock}}}, \mathbf{C})$  to use in the estimators. Similar to the experiments detailed in [Póczos et al. \(2012\)](#), we construct  $\mathbf{Y}^{\text{ref}}$  with a comparable sample size as  $\mathbf{X}^{\text{mock}}$ : 2000 and 10,000 for the  $P_\ell$  and  $\zeta$  analyses respectively.

In Figure 2 we compare the distribution of Rényi- $\alpha$  (left) and  $KL$  (right) divergence estimates (orange)  $\hat{D}_{R\alpha}$  and  $\hat{D}_{KL}$  between the mock data  $\mathbf{X}^{\text{mock}}$  and a reference sample  $\mathbf{Y}^{\text{ref}}$  for the  $P_\ell(k)$  (top) and  $\zeta(N)$  (bottom) analyses. As a reference point for the comparison, we include in Figure 2 (blue) the distribution of Rényi- $\alpha$  and  $KL$  divergence estimates if  $\mathbf{X}^{\text{mock}}$  were actually sampled from  $\mathcal{N}(\overline{\mathbf{X}^{\text{mock}}}, \mathbf{C})$  and  $\mathbf{Y}^{\text{ref}}$ . All distributions were constructed using 100 divergence estimates.

The discrepancy between the blue and orange distributions illustrates the non-Gaussianity of  $p(x)$  sampled by  $\mathbf{X}^{\text{mock}}$ .

- Describe the discrepancy.

## 5. A More Accurate Likelihood

- Gaussian Mixture Model

- expectation maximization algorithm [Dempster et al. \(1977\)](#)
- Bayesian Information Criteria
- Figure illustrating both methods on highest N GMF bin

### 5.1. Independent Component Analysis

Curse of dimensionality! 2048 mocks in Beutler not enough to directly estimate the 37-dimensional space, so we use independent component analysis [Hartlap et al. \(2009\)](#)

- equation
- Similar figure to [Hartlap et al. \(2009\)](#) that tests the independence?
- Kernel Density Estimation
- 

Figure [2](#)

## 6. Impact on Parameter Inference

### 6.1. MCMC

- details of each of the MCMC runs

### 6.2. Importance Sampling

- equations explaining importance sampling framework

Figure [3](#)

Figure [4](#)

## 7. Discussion

- Will it matter for future surveys?
- Likelihood free inference (cite justin's paper)

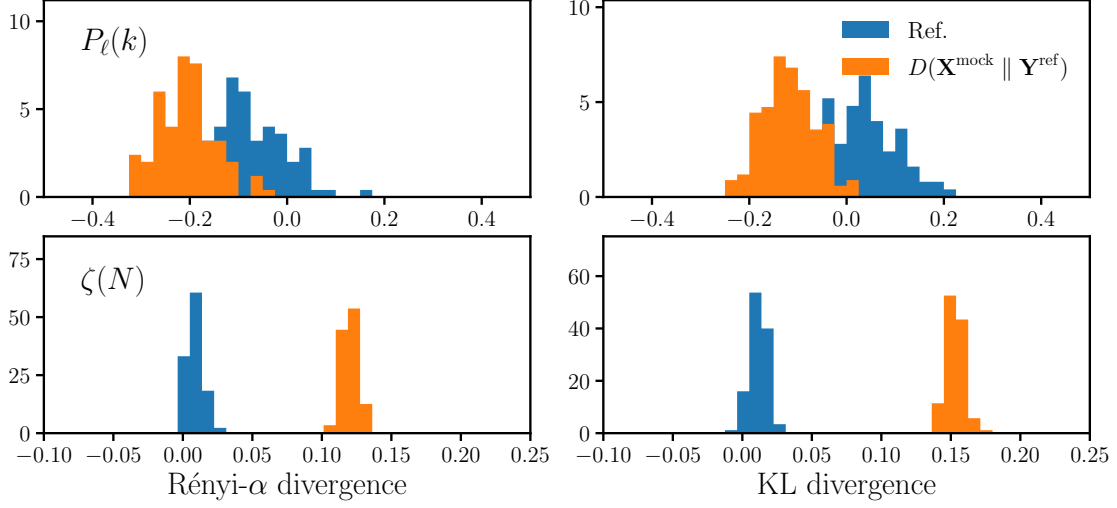


Fig. 1.— Rényi- $\alpha$  and KL divergence estimates, ( $D_{R\alpha}$  and  $D_{KL}$ ), between the mock data  $\mathbf{X}^{\text{mock}}$  and a reference sample  $\mathbf{Y}^{\text{ref}}$  for the  $P_\ell(k)$  (left) and  $\zeta(N)$  (right) analyses.

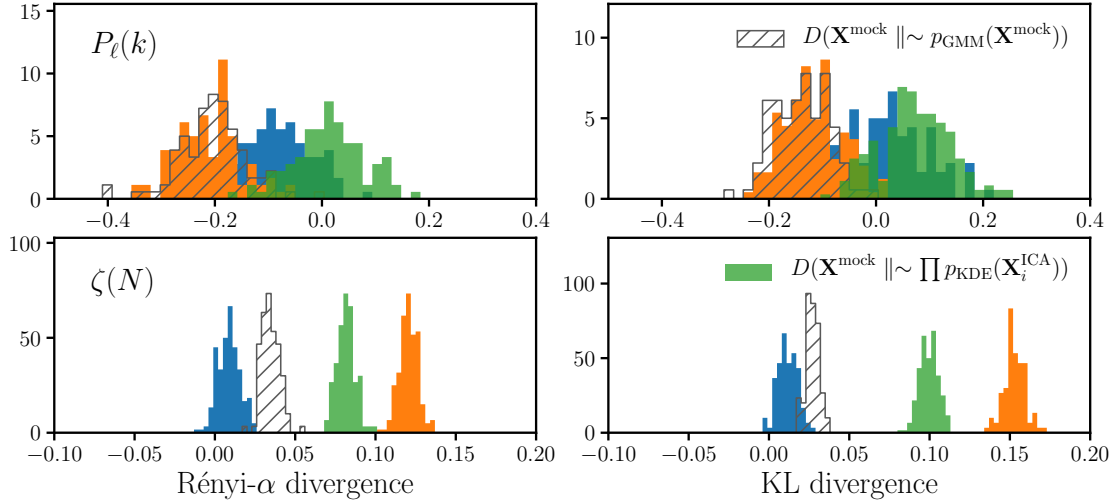


Fig. 2.—

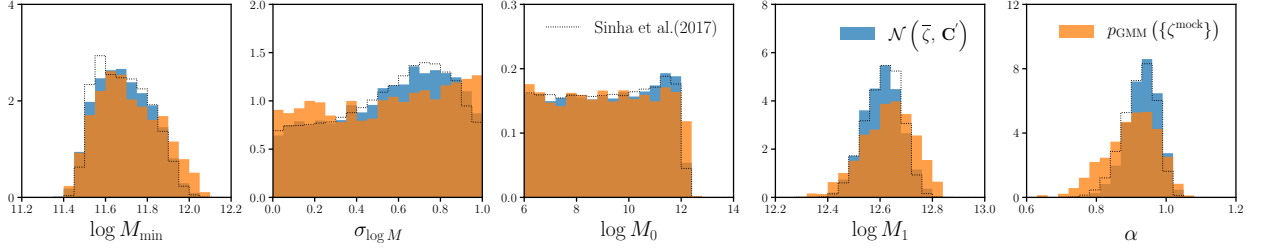


Fig. 3.—

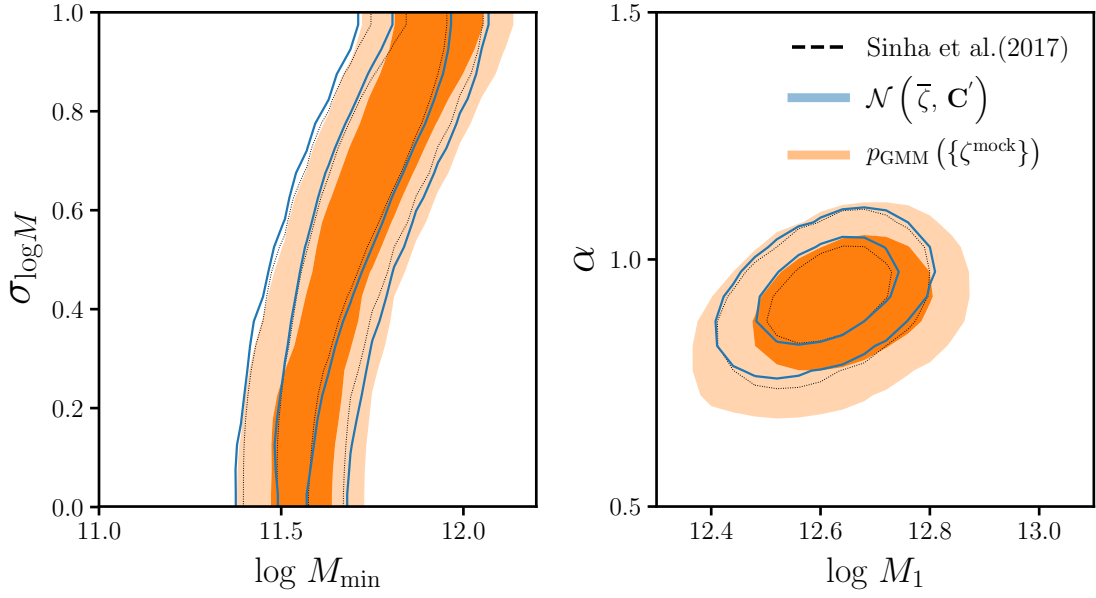


Fig. 4.—

## 8. Summary

### Acknowledgements

It’s a pleasure to thank Simone Ferraro, David W. Hogg, Emmaneul Schaan, Roman Scocimarro Zachary Slepian

### REFERENCES

- Alam, S., Ata, M., Bailey, S., et al. 2017, [Monthly Notices of the Royal Astronomical Society](#), 470, 2617
- Beutler, F., Seo, H.-J., Saito, S., et al. 2017, [Monthly Notices of the Royal Astronomical Society](#), 466, 2242
- Dempster, A. P., Laird, N. M., & Rubin, D. B. 1977, *Journal of the Royal Statistical Society. Series B (Methodological)*, 39, 1
- Grieb, J. N., Sánchez, A. G., Salazar-Albornoz, S., et al. 2017, [Monthly Notices of the Royal Astronomical Society](#), 467, 2085
- Guo, H., Zehavi, I., & Zheng, Z. 2012, [The Astrophysical Journal](#), 756, 127
- Hahn, C., Scocimarro, R., Blanton, M. R., Tinker, J. L., & Rodríguez-Torres, S. A. 2017, [Monthly Notices of the Royal Astronomical Society](#), 467, 1940
- Hartlap, J., Schrabback, T., Simon, P., & Schneider, P. 2009, [Astronomy and Astrophysics](#), 504, 689
- Kazin, E. A., Koda, J., Blake, C., et al. 2014, [Monthly Notices of the Royal Astronomical Society](#), 441, 3524
- Kitaura, F.-S., Gil-Marín, H., Scóccola, C. G., et al. 2015, [Monthly Notices of the Royal Astronomical Society](#), 450, 1836
- Kitaura, F.-S., & Heß, S. 2013, [Monthly Notices of the Royal Astronomical Society](#), 435, L78
- Kitaura, F.-S., Yepes, G., & Prada, F. 2014, [Monthly Notices of the Royal Astronomical Society](#), 439, L21
- Kitaura, F.-S., Rodríguez-Torres, S., Chuang, C.-H., et al. 2016, [Monthly Notices of the Royal Astronomical Society](#), 456, 4156
- Klypin, A., Yepes, G., Gottlöber, S., Prada, F., & Heß, S. 2016, [Monthly Notices of the Royal Astronomical Society](#), 457, 4340
- Krishnamurthy, A., Kandasamy, K., Poczos, B., & Wasserman, L. 2014, arXiv:1402.2966 [math, stat], [arXiv:1402.2966 \[math, stat\]](#)
- Mohammed, I., Seljak, U., & Vlah, Z. 2017, [Monthly Notices of the Royal Astronomical Society](#), 466, 780
- Norberg, P., Baugh, C. M., Gaztañaga, E., & Croton, D. J. 2009, [Monthly Notices of the Royal Astronomical Society](#), 396, 19



- Ntampaka, M., Trac, H., Sutherland, D. J., et al. 2015, [The Astrophysical Journal](#), 803, 50
- . 2016, [The Astrophysical Journal](#), 831, 135
- Parkinson, D., Riemer-Sørensen, S., Blake, C., et al. 2012, [Physical Review D](#), 86, 103518
- Pinol, L., Cahn, R. N., Hand, N., Seljak, U., & White, M. 2017, [Journal of Cosmology and Astroparticle Physics](#), 2017, 008
- Póczos, B., Xiong, L., Sutherland, D. J., & Schneider, J. 2012, in [2012 IEEE Conference on Computer Vision and Pattern Recognition](#), 2989
- Rodríguez-Torres, S. A., Chuang, C.-H., Prada, F., et al. 2016, [Monthly Notices of the Royal Astronomical Society](#), 460, 1173
- Ross, A. J., Beutler, F., Chuang, C.-H., et al. 2017, [Monthly Notices of the Royal Astronomical Society](#), 464, 1168
- Sellentin, E., Jaffe, A. H., & Heavens, A. F. 2017, arXiv:1709.03452 [astro-ph, stat], [arXiv:1709.03452 \[astro-ph, stat\]](#)
- Sinha, M., Berlind, A. A., McBride, C. K., et al. 2017, arXiv:1708.04892 [astro-ph], [arXiv:1708.04892 \[astro-ph\]](#)
- Slepian, Z., Eisenstein, D. J., Brownstein, J. R., et al. 2017, [Monthly Notices of the Royal Astronomical Society](#), 469, 1738
- Tinker, J. L., & et al. in preparation
- Vargas-Magaña, M., Ho, S., Xu, X., et al. 2014, [Monthly Notices of the Royal Astronomical Society](#), 445, 2
- Wang, Q., Sanjeev, K., & Sergio, V. 2009, [IEEE TRANSACTIONS ON INFORMATION THEORY](#), 55, 2392
- Zhao, C., Kitaura, F.-S., Chuang, C.-H., et al. 2015, [Monthly Notices of the Royal Astronomical Society](#), 451, 4266