

What can random matrices tell us about algorithms?

Case studies.

Tom Trogdon

trogdon@uw.edu

University of Washington

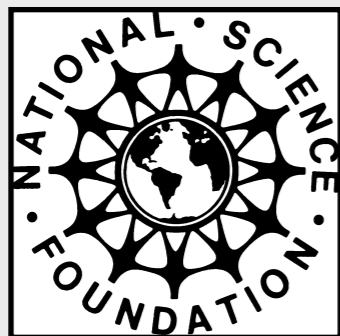
Department of Applied Mathematics



Acknowledgements

Joint work with:

- Percy Deift (NYU)
- Xiukai Ding (UC Davis)
- Tyler Chen (NYU)
- Elliot Paquette (McGill)



NSF Funding via:

CAREER: Numerical linear algebra, random matrix theory and applications

Perspectives on *Matrix Computations*:
Theoretical Computer Science Meets
Numerical analysis

Perspectives on *Matrix Computations*:
Theoretical Computer Science Meets
Numerical analysis

... Meets *Mathematical Physics*?

Perspectives on Matrix Computations: Theoretical Computer Science Meets Numerical analysis

... Meets Mathematical Physics?

P Deift, T Nanda, and C Tomei. Ordinary differential equations and the symmetric eigenvalue problem. SIAM J. on Numer. Anal., 20:1–22, 1983

Plan

- Lanczos/CG and the power method on sample covariance matrices
- Lanczos and finite-precision effects
- Riemann-Hilbert analysis for CG/Lanczos*



Building blocks of random matrix theory

The real Ginibre Ensemble, $\text{Gin}_{\mathbb{R}}(N, M)$, is the matrix

$$Y = \left(\frac{Y_{ij}}{\sqrt{M}} \right)_{1 \leq i \leq N, 1 \leq j \leq M}, \quad Y_{ij} \text{ iid standard normal random variables.}$$

The complex Ginibre Ensemble, $\text{Gin}_{\mathbb{C}}(N, M)$, is the matrix

$$X = \frac{1}{\sqrt{2}} (Y_1 + iY_2), \quad Y_1, Y_2 \text{ independent } \text{Gin}_{\mathbb{R}}(N, M).$$

If $N = M$ we use $\text{Gin}_{\mathbb{F}}(N)$.

The Gaussian Unitary Ensemble, $\text{GUE}(N)$, is the matrix

$$H = \frac{1}{\sqrt{2}} (X + X^*), \quad X \sim \text{Gin}_{\mathbb{C}}(N).$$

Sample covariance matrices (SCMs)

The classical examples of sample covariance matrices are:

The real Wishart Ensemble, $\text{Wishart}_{\mathbb{R}}(N, M)$:

$$W = XX^T, \quad X \sim \text{Gin}_{\mathbb{R}}(N, M).$$

The complex Wishart Ensemble, $\text{Wishart}_{\mathbb{C}}(N, M)$:

$$W = XX^*, \quad X \sim \text{Gin}_{\mathbb{C}}(N, M).$$

Sample covariance matrices (SCMs)

The classical examples of sample covariance matrices are:

The real Wishart Ensemble, $\text{Wishart}_{\mathbb{R}}(N, M)$:

$$W = XX^T, \quad X \sim \text{Gin}_{\mathbb{R}}(N, M).$$

The complex Wishart Ensemble, $\text{Wishart}_{\mathbb{C}}(N, M)$:

$$W = XX^*, \quad X \sim \text{Gin}_{\mathbb{C}}(N, M).$$

But there is, in principle, no reason why the entries of X should be Gaussian, and no reason why the entries in each column need to be independent.

Example: X could be an $N \times M$ matrix of iid random variables, $X_{11} = \pm 1/\sqrt{M}$ with equal probability.

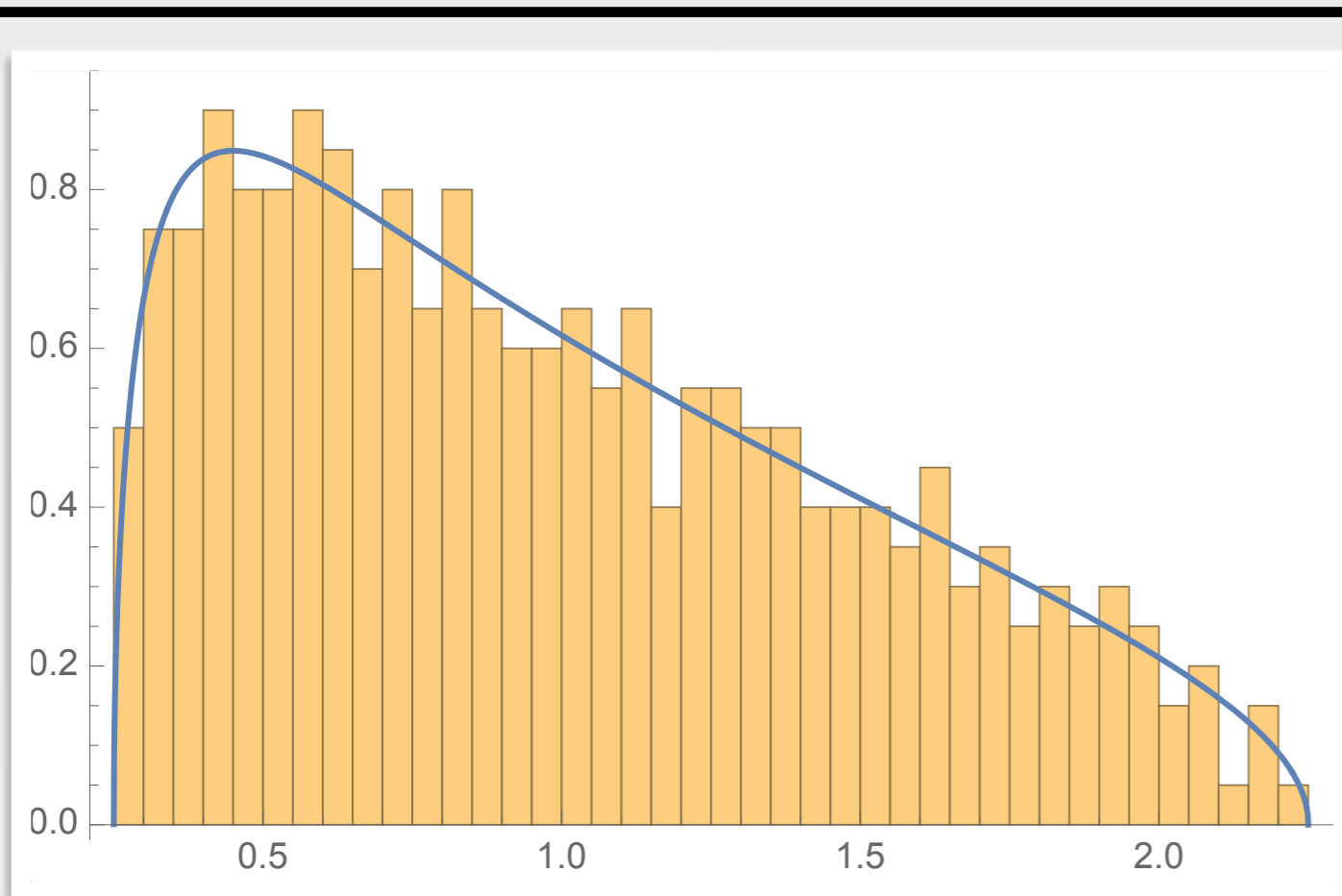
Sample covariance matrices (SCMs)

The classical examples of sample covariance matrices are:

The real Wishart Ensemble, $\text{Wishart}_{\mathbb{R}}(N, M)$:

$$W = XX^T, \quad X \sim \text{Gin}_{\mathbb{R}}(N, M).$$

The complex Wishart Ensemble, $\text{Wishart}_{\mathbb{C}}(N, M)$:



The Marchenko–Pastur distribution gives the macroscopic behavior of the spectrum:

$$p_d(x) = \frac{1}{2\pi d} \frac{\sqrt{[(\lambda_+ - x)(x - \lambda_-)]_+}}{|x|}$$

$$\lambda_{\pm} = (1 \pm \sqrt{d})^2$$

$$M \sim N/d$$

$$\frac{1}{N} \sum_{j=1}^N \delta_{\lambda_j} \rightarrow p_d(x) dx$$

CG and Wishart: An exactly solvable case



The CG algorithm and orthogonal polynomials

The CG algorithm applied to $W\mathbf{x} = \mathbf{b}$ is equivalent to the polynomial minimization problem

$$\min_{P: \text{degree } k, P(0)=1} \|W^{-1}P(W)\mathbf{b}\|_W^2 = \min_{P: \text{degree } k, P(0)=1} \int \frac{P(\lambda)^2}{\lambda} \mu_{W,\mathbf{b}}(d\lambda),$$

where

$$\mu_{W,\mathbf{b}} = \sum_{j=1}^N |\langle \mathbf{q}_j, \mathbf{b} \rangle|^2 \delta_{\lambda_j},$$

is the so-called eigenvector spectral distribution (VESD).

M Hestenes and E Steifel. Method of Conjugate Gradients for Solving Linear Systems. J. Research Nat. Bur. Standards, 20:409–436, 1952

Z Bai and J W Silverstein. Spectral Analysis of Large Dimensional Random Matrices. Springer Series in Statistics. Springer New York, New York, NY, 2010



The CG algorithm and orthogonal polynomials

The CG algorithm applied to $W\mathbf{x} = \mathbf{b}$ is equivalent to the polynomial minimization problem

$$\min_{P: \text{degree } k, P(0)=1} \|W^{-1}P(W)\mathbf{b}\|_W^2 = \min_{P: \text{degree } k, P(0)=1} \int \frac{P(\lambda)^2}{\lambda} \mu_{W,\mathbf{b}}(d\lambda),$$

where

$$\mu_{W,\mathbf{b}} = \sum_{j=1}^N |\langle \mathbf{q}_j, \mathbf{b} \rangle|^2 \delta_{\lambda_j},$$

is the so-called eigenvector spectral distribution (VESD).

M Hestenes and E Steifel. Method of Conjugate Gradients for Solving Linear Systems. J. Research Nat. Bur. Standards, 20:409–436, 1952

Z Bai and J W Silverstein. Spectral Analysis of Large Dimensional Random Matrices. Springer Series in Statistics. Springer New York, New York, NY, 2010

A useful (Riemann–Hilbert) tool for orthogonal polynomials:

A S Fokas, A R Its, and A V Kitaev. The isomonodromy approach to matrix models in 2D quantum gravity. Comm. Math. Phys., 147(2):395–430, 1992

P Deift. Orthogonal Polynomials and Random Matrices: a Riemann-Hilbert Approach. Amer. Math. Soc., Providence, RI, 2000



Beyond Wishart: Universal fluctuations

$$\|\mathbf{r}_k(W, \mathbf{b})\|_2^2 \sim \left(\frac{N}{M}\right)^k \left(1 + \sqrt{\frac{2k}{\beta M}} \sqrt{1 + \frac{N}{M}} \mathcal{N}(0, 1)\right)$$

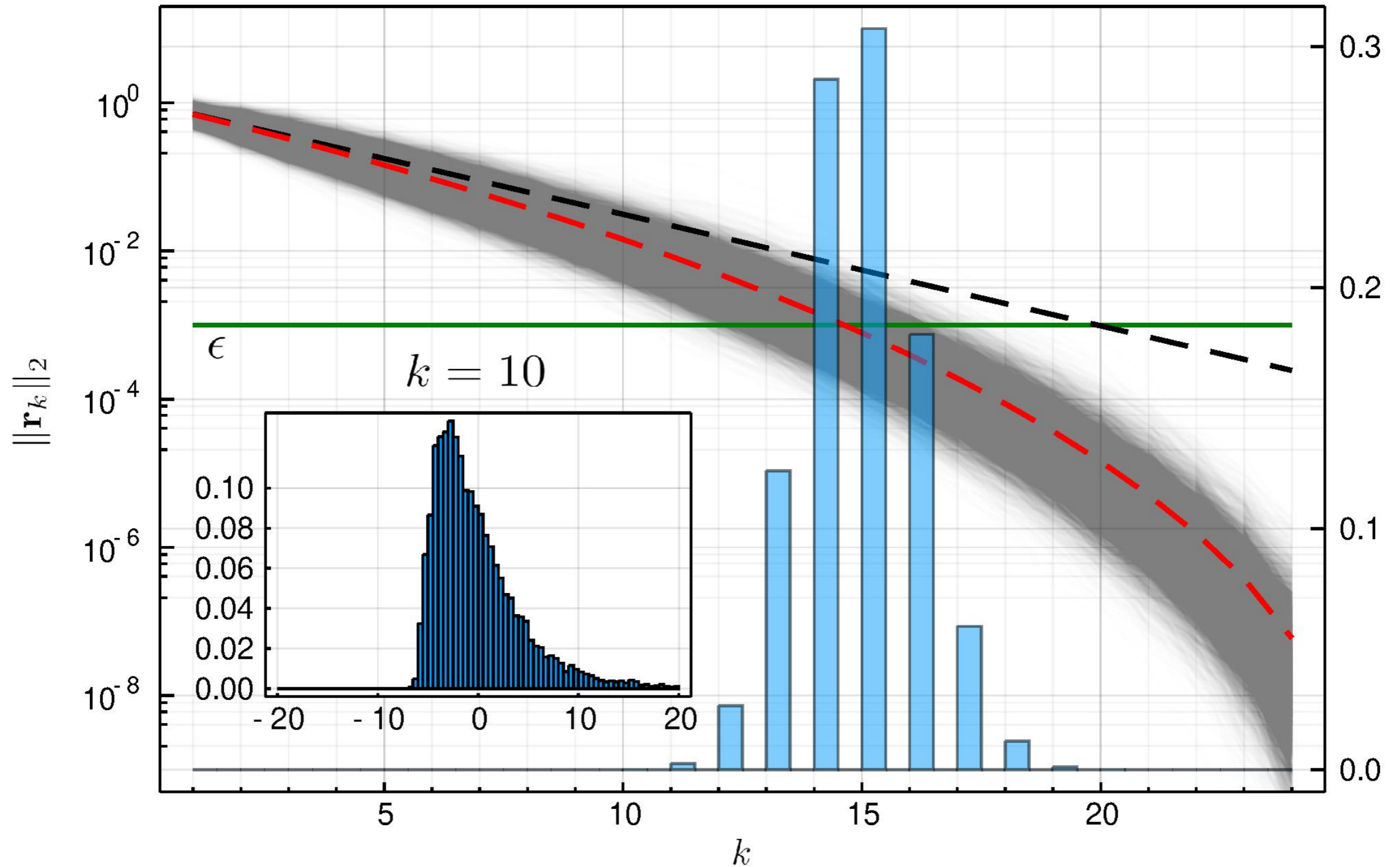
We will also consider the halting time:

$$T(W, \mathbf{b}, \epsilon) := \min\{k : \|\mathbf{r}_k(W, \mathbf{b})\| < \epsilon\}.$$



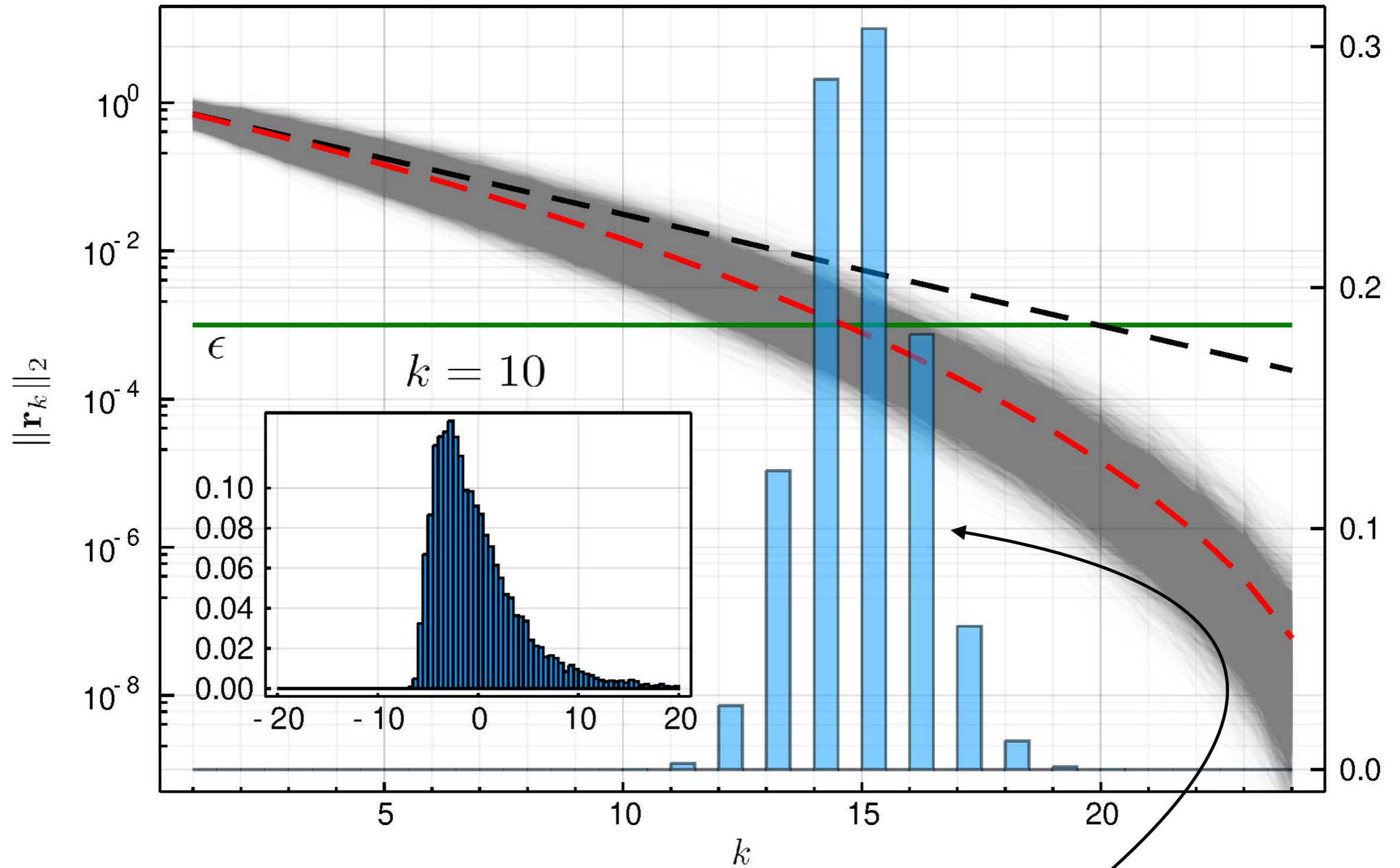
Sample covariance with trivial covariance

$$N = 25, M = 50, d = 0.50$$



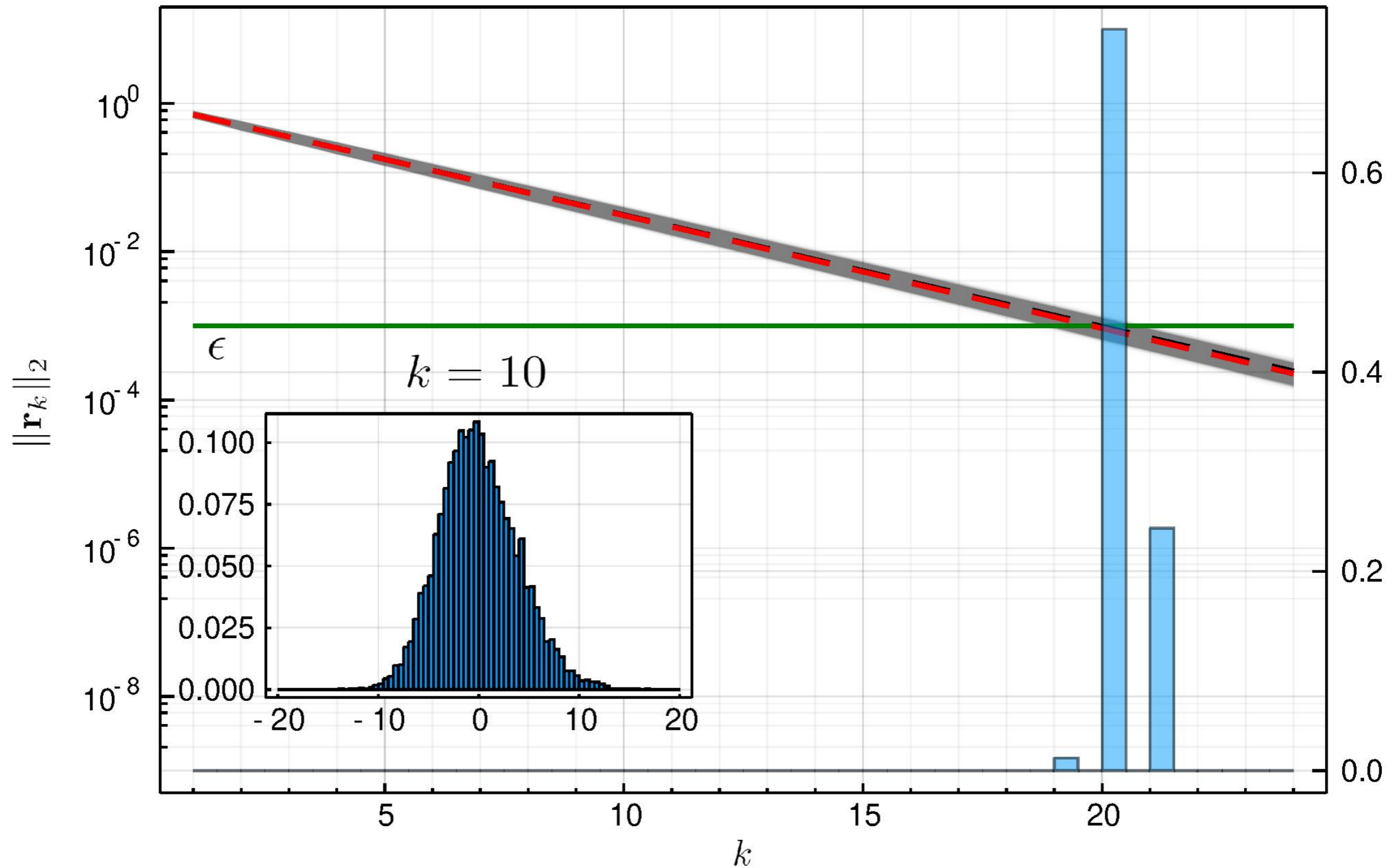
Sample covariance with trivial covariance

$$N = 25, M = 50, d = 0.50$$



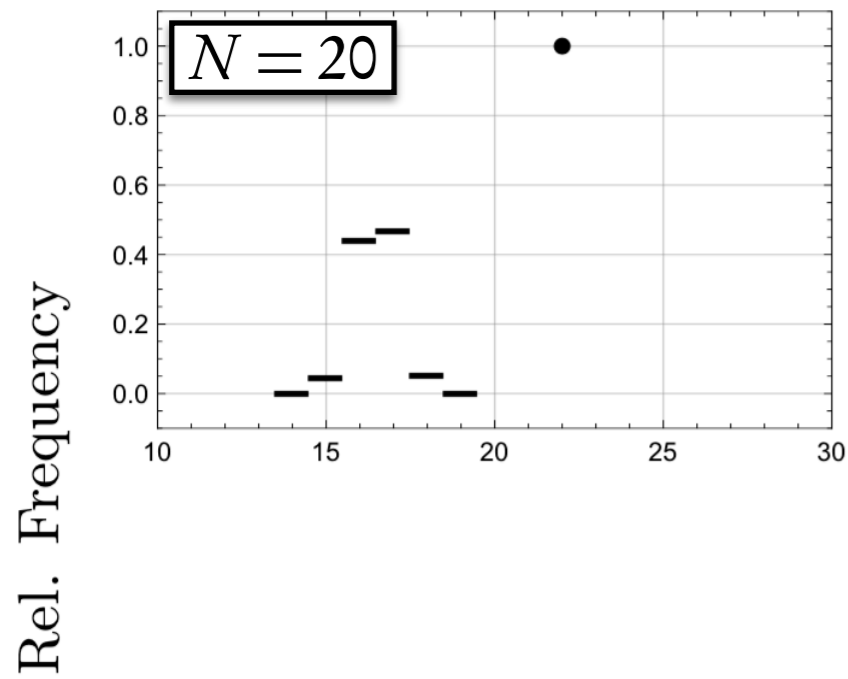
Sample covariance with trivial covariance

$N = 1000, M = 2000, d = 0.50$



Concentration occurs if $\frac{\log \epsilon^{-1}}{\sqrt{M}} \ll 1$





Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

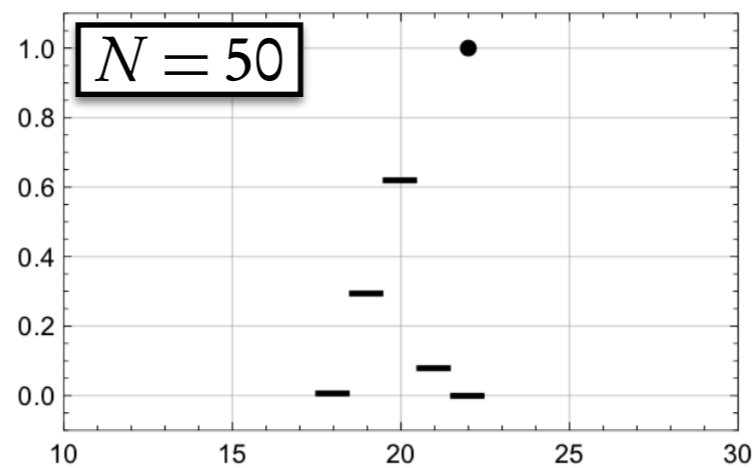
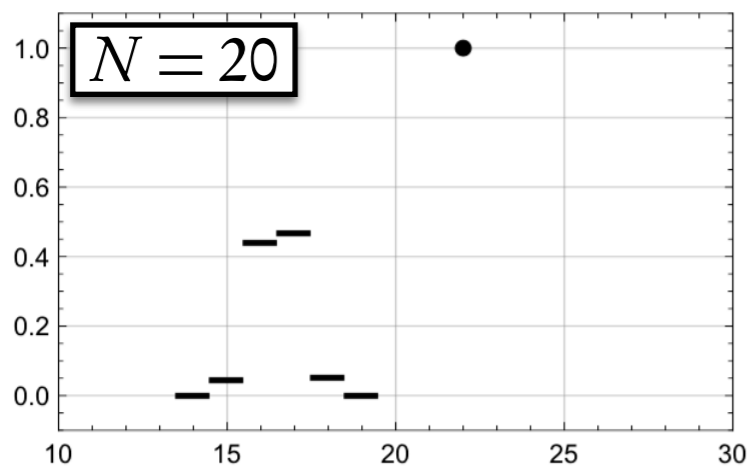
Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$



Rel. Frequency



Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

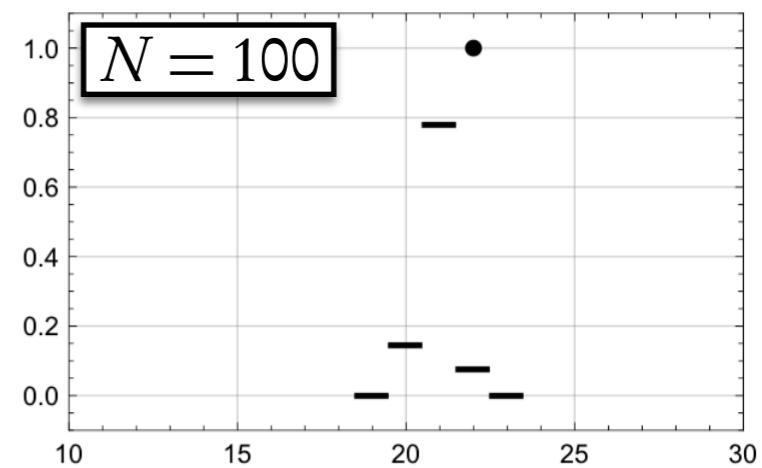
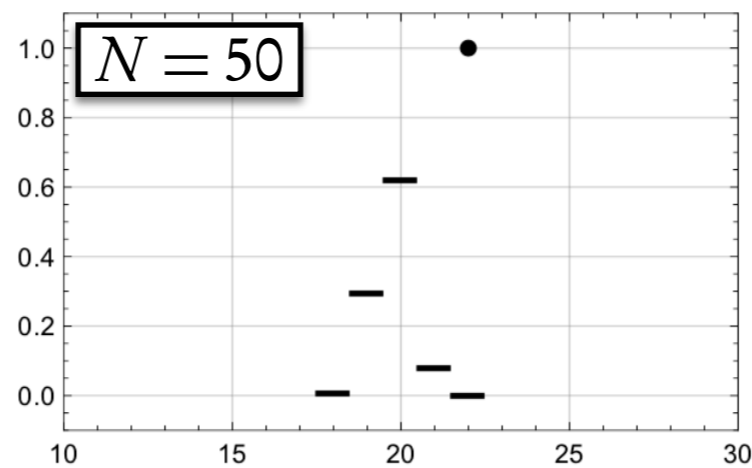
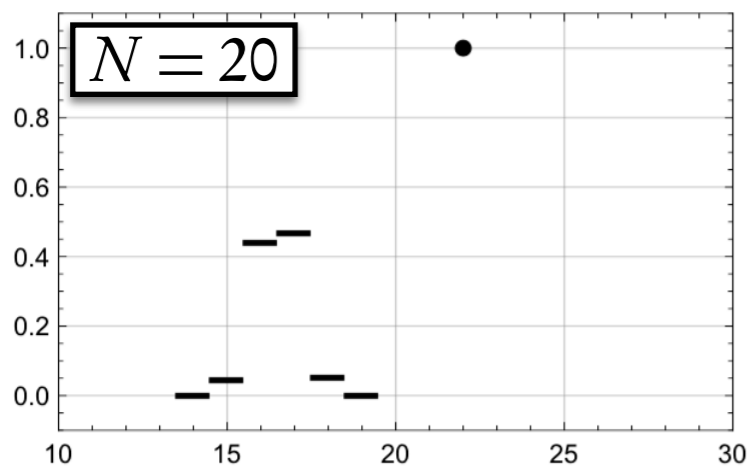
Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$



Rel. Frequency



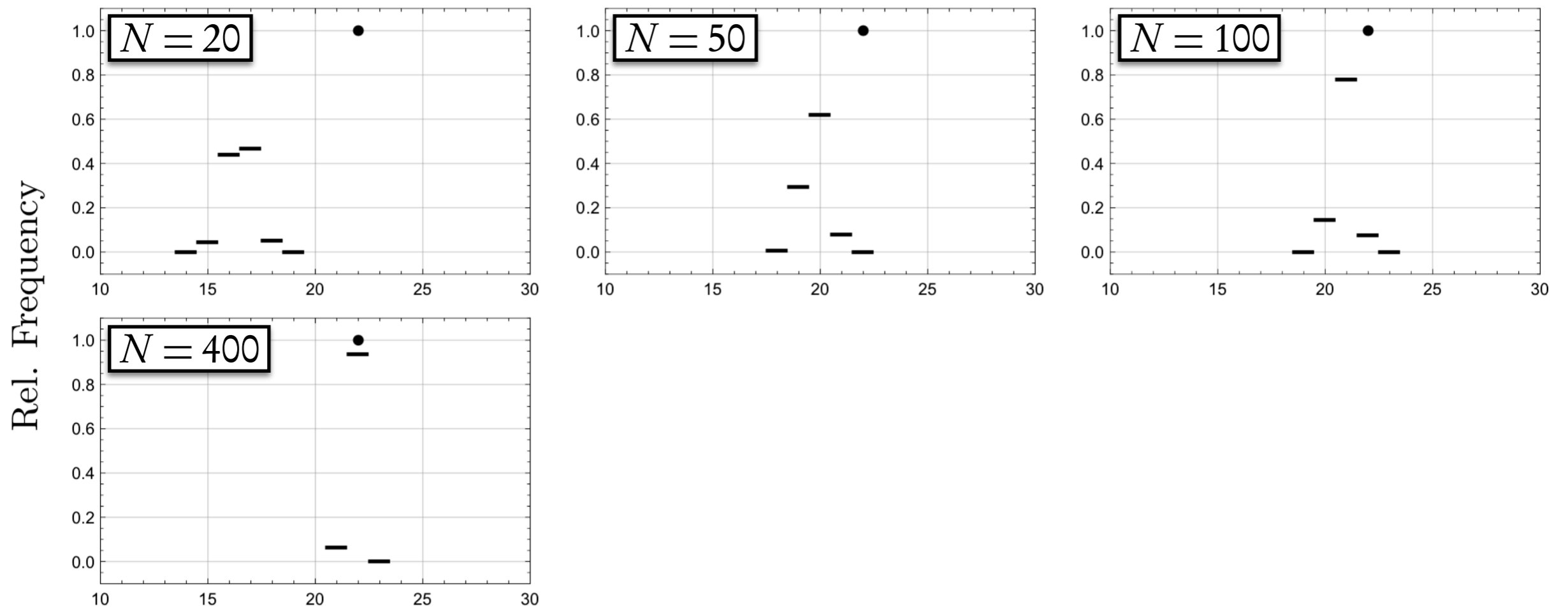
Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$





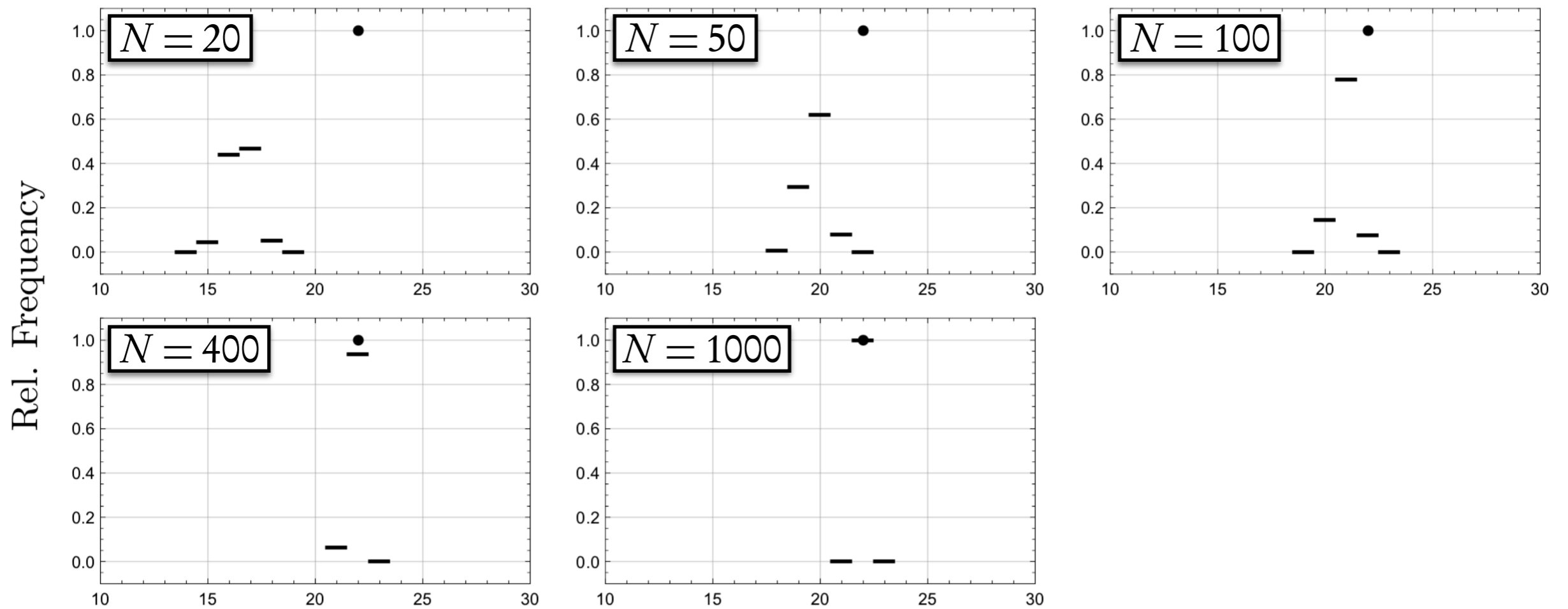
Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$





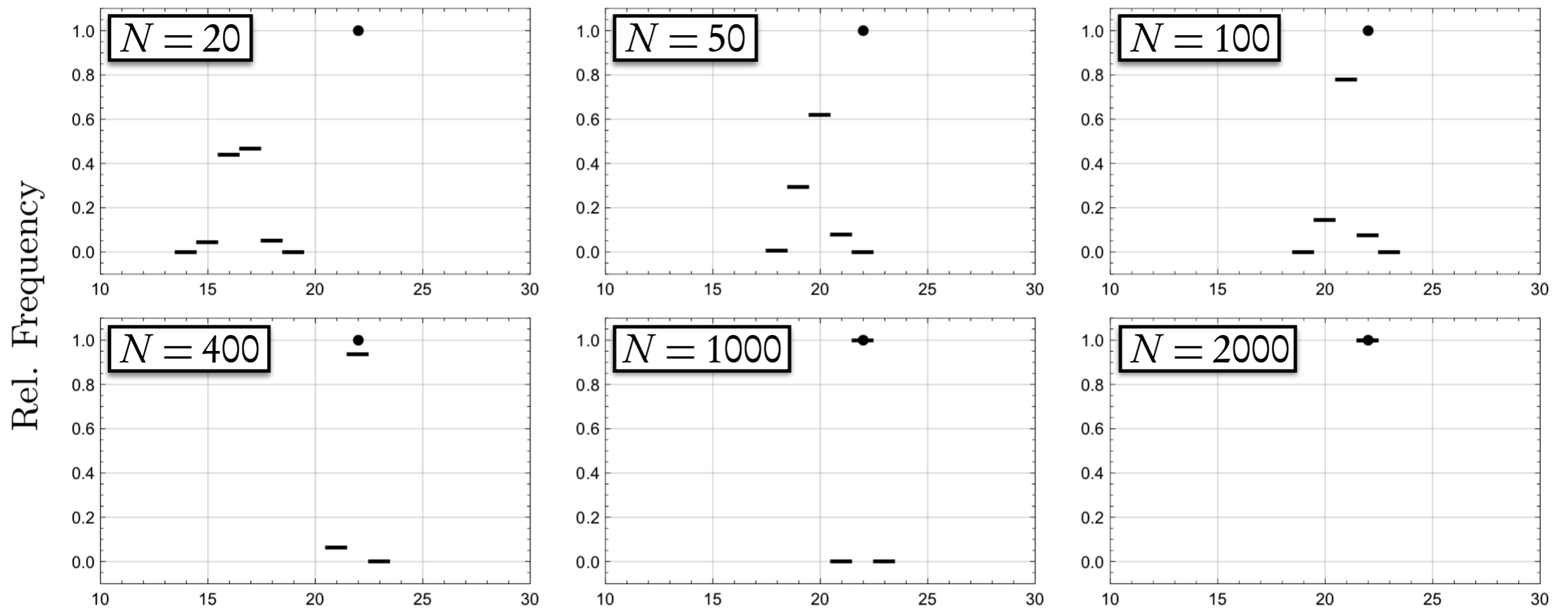
Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$





Statistics for $T(W, \mathbf{b}, \epsilon) = \min\{k : \|\mathbf{r}_k\|_2 < \epsilon\}$

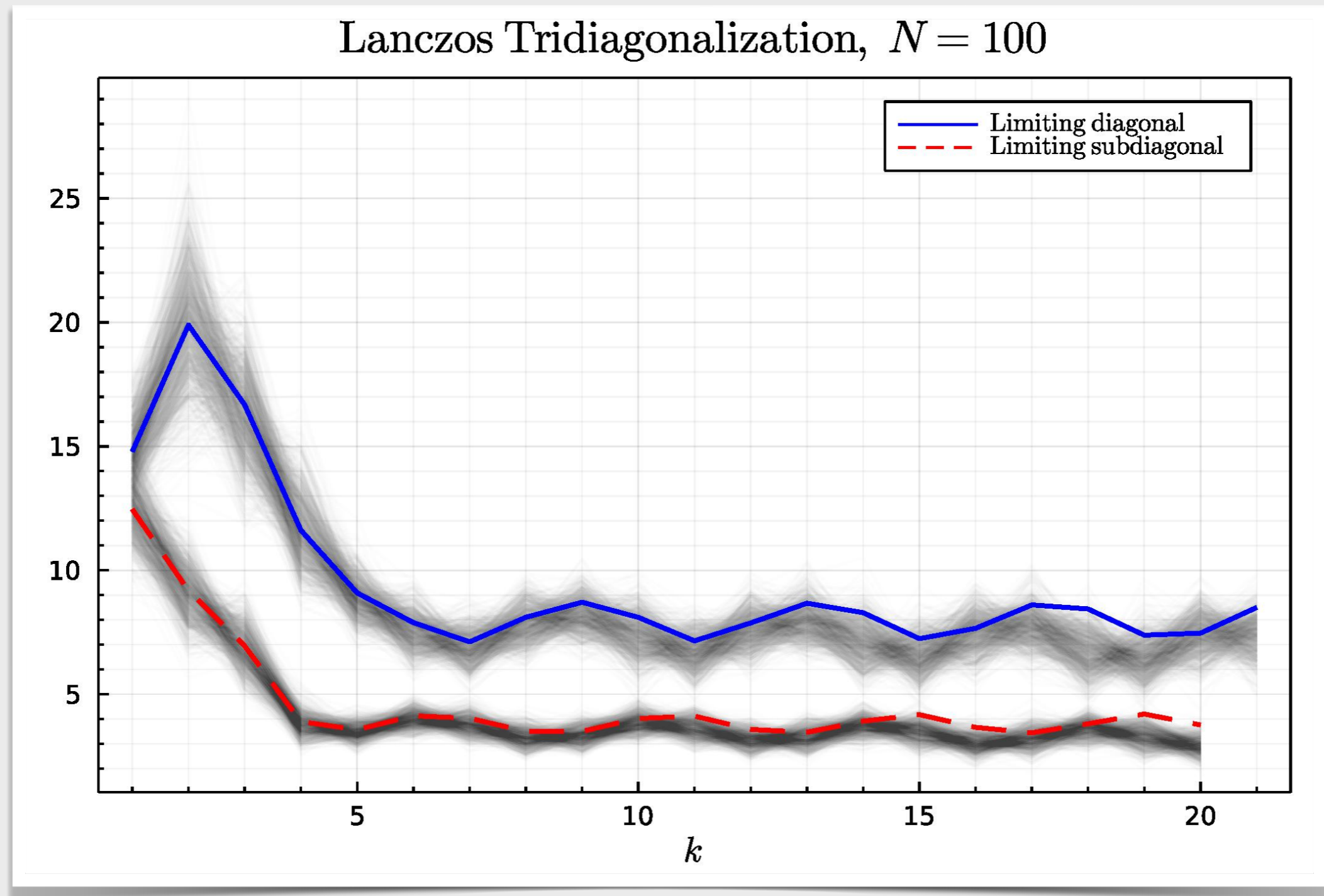
Wishart $_{\mathbb{R}}(N, M)$, $d = 0.2$, $\epsilon \approx 6 \times 10^{-8}$
20,000 samples

$$\forall k \quad \epsilon \neq d^{k/2} \quad \mathbb{P}\left(T(W, \mathbf{b}, \epsilon) \neq \left\lceil \frac{2 \log \epsilon}{\log d} \right\rceil\right) \leq C e^{-cN}$$

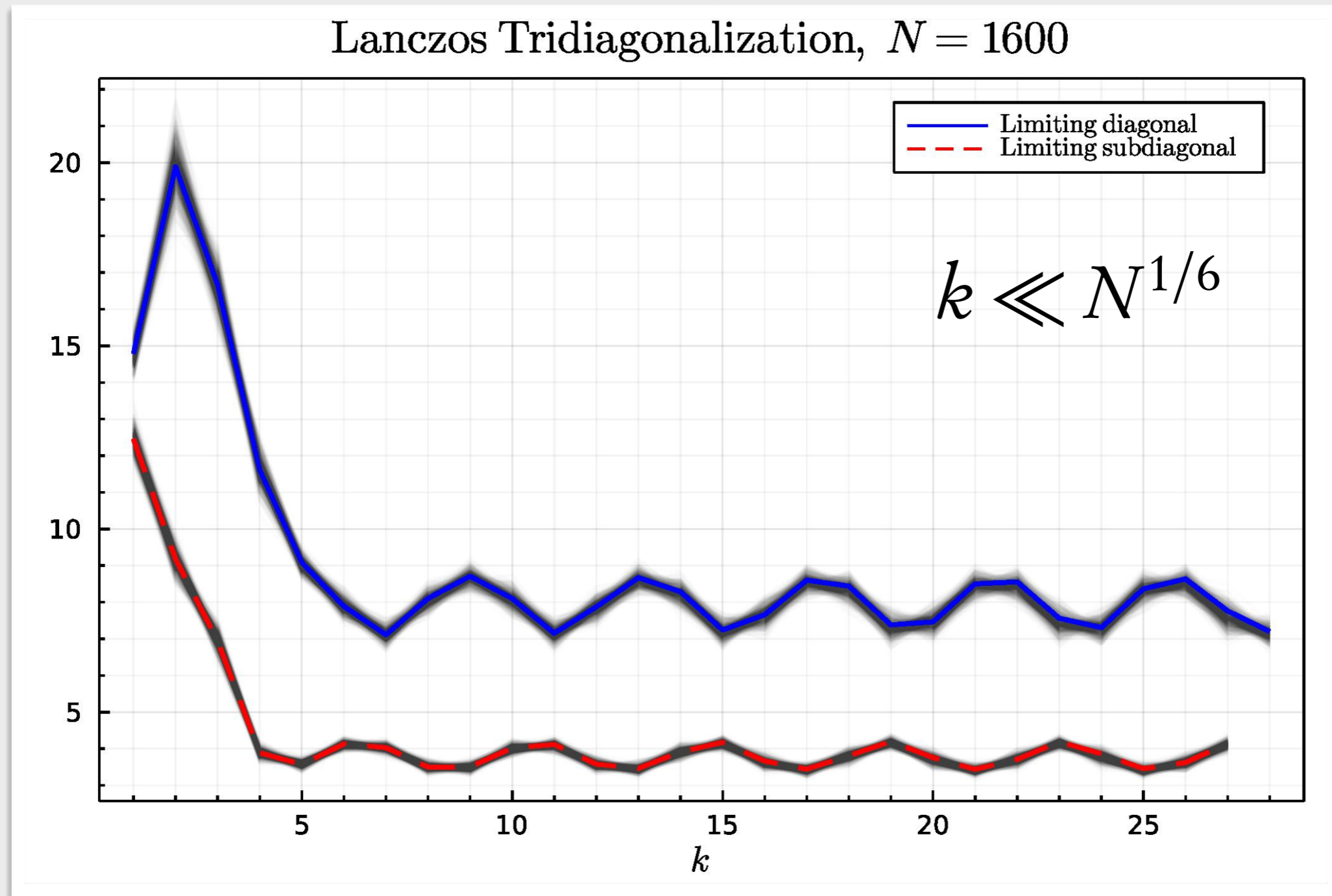
$$\epsilon = d^{k/2} \quad \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k) = \frac{1}{2} = \lim_{N \rightarrow \infty} \mathbb{P}(T(W, \mathbf{b}, \epsilon) = k + 1)$$



Spiked covariance with multiple bulk components



Spiked covariance with multiple bulk components



X Ding and T T. The conjugate gradient algorithm on a general class of spiked covariance matrices. *Quarterly of Applied Mathematics*, 80(1):99–155, nov 2021

X Ding and T T. A Riemann–Hilbert approach to the perturbation theory for orthogonal polynomials: Applications to numerical linear algebra and random matrix theory. *arXiv preprint 2112.12354*, pages 1–77, 2021



The power method

$$\nu_k(W, \mathbf{b}) = \frac{\int \lambda^{2k-1} \mu_{W, \mathbf{b}}(d\lambda)}{\int \lambda^{2k-2} \mu_{W, \mathbf{b}}(d\lambda)}, \quad k = 1, 2, \dots$$

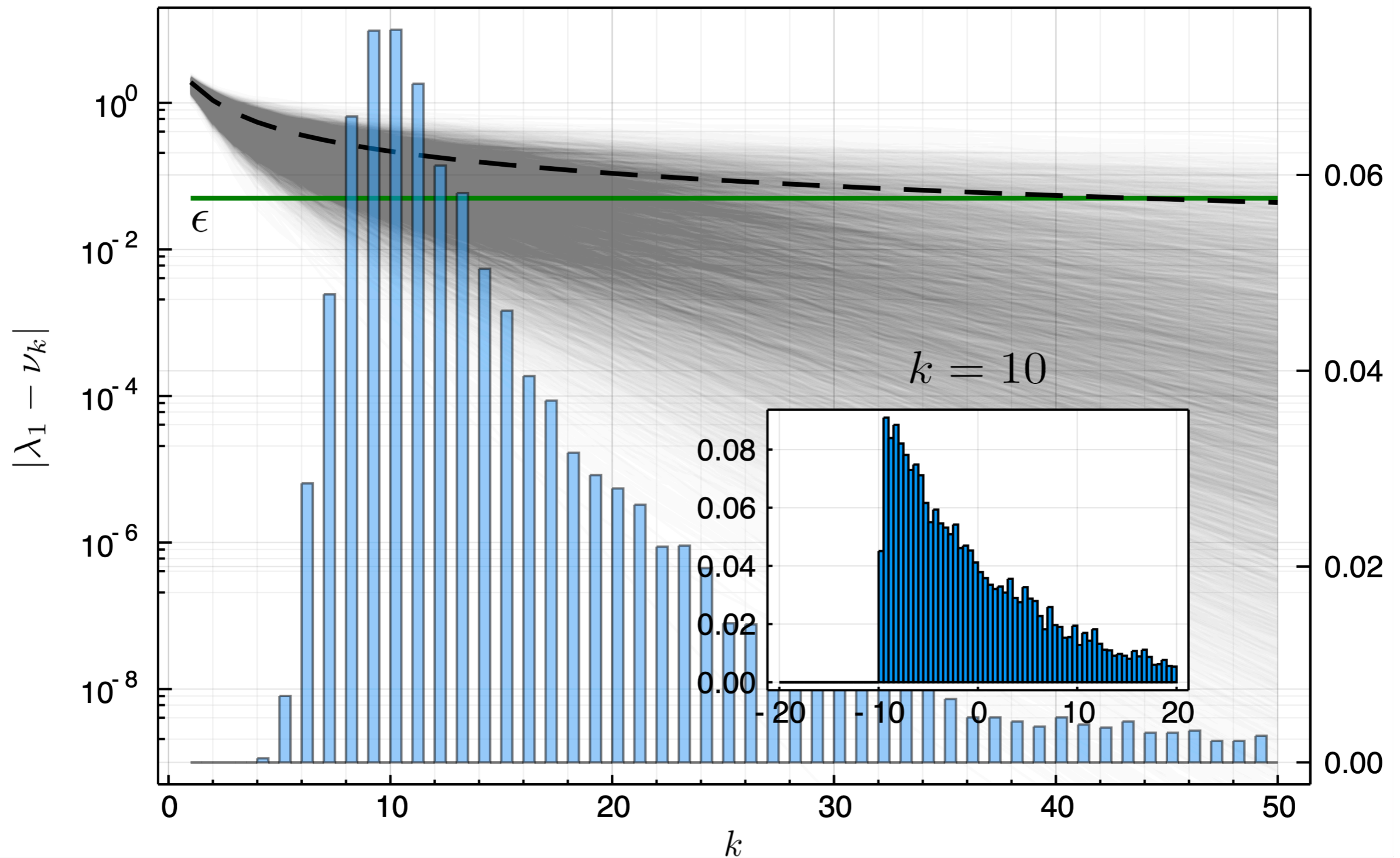
$$\nu_k(W, \mathbf{b}) \xrightarrow{N \rightarrow \infty} \frac{\int \lambda^{2k-1} \mu_{\text{MP}}(d\lambda)}{\int \lambda^{2k-2} \mu_{\text{MP}}(d\lambda)}, \quad k = 1, 2, \dots$$

Consider the measure of error

$$E_k(W, \mathbf{b}) = |\lambda_{\max}(W) - \nu_k(W, \mathbf{b})|, \quad k = 1, 2, \dots$$



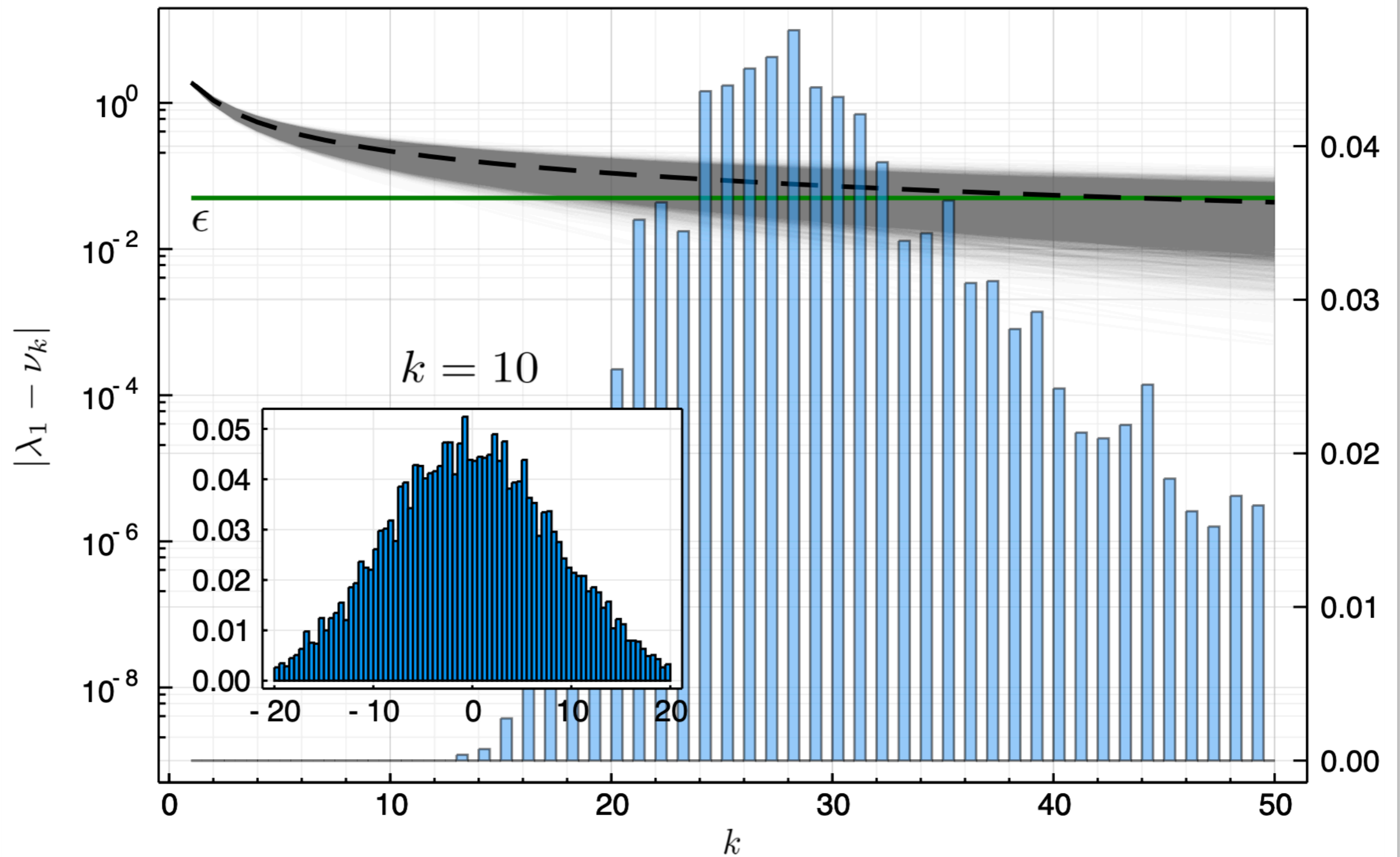
$$N = 50, M = 100, d = 0.50$$



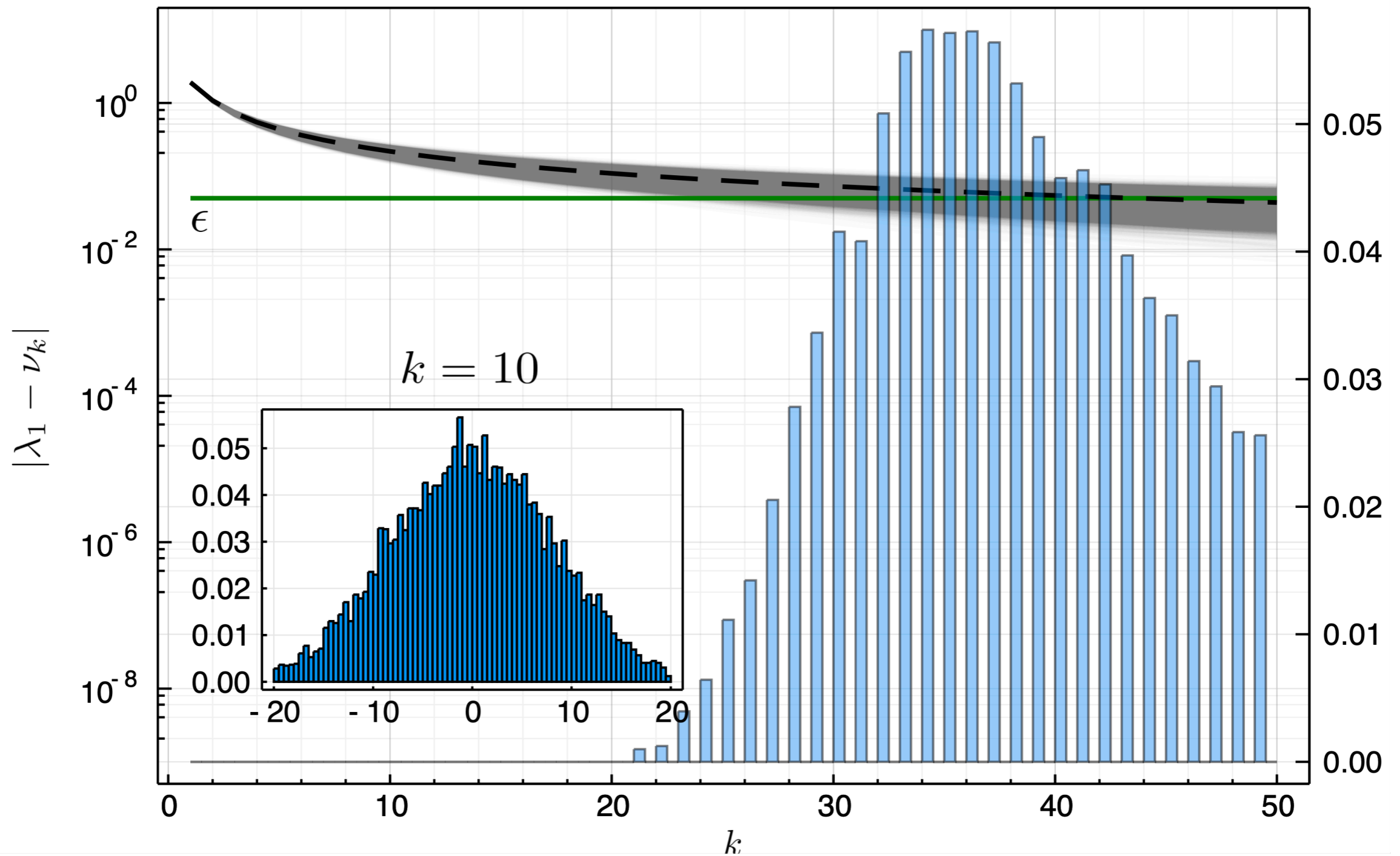
J Kuczyński and H Woźniakowski. Estimating the Largest Eigenvalue by the Power and Lanczos Algorithms with a Random Start. *SIAM Journal on Matrix Analysis and Applications*, 13(4):1094–1122, oct 1992



$N = 1000, M = 2000, d = 0.50$



$N = 4000, M = 8000, d = 0.50$

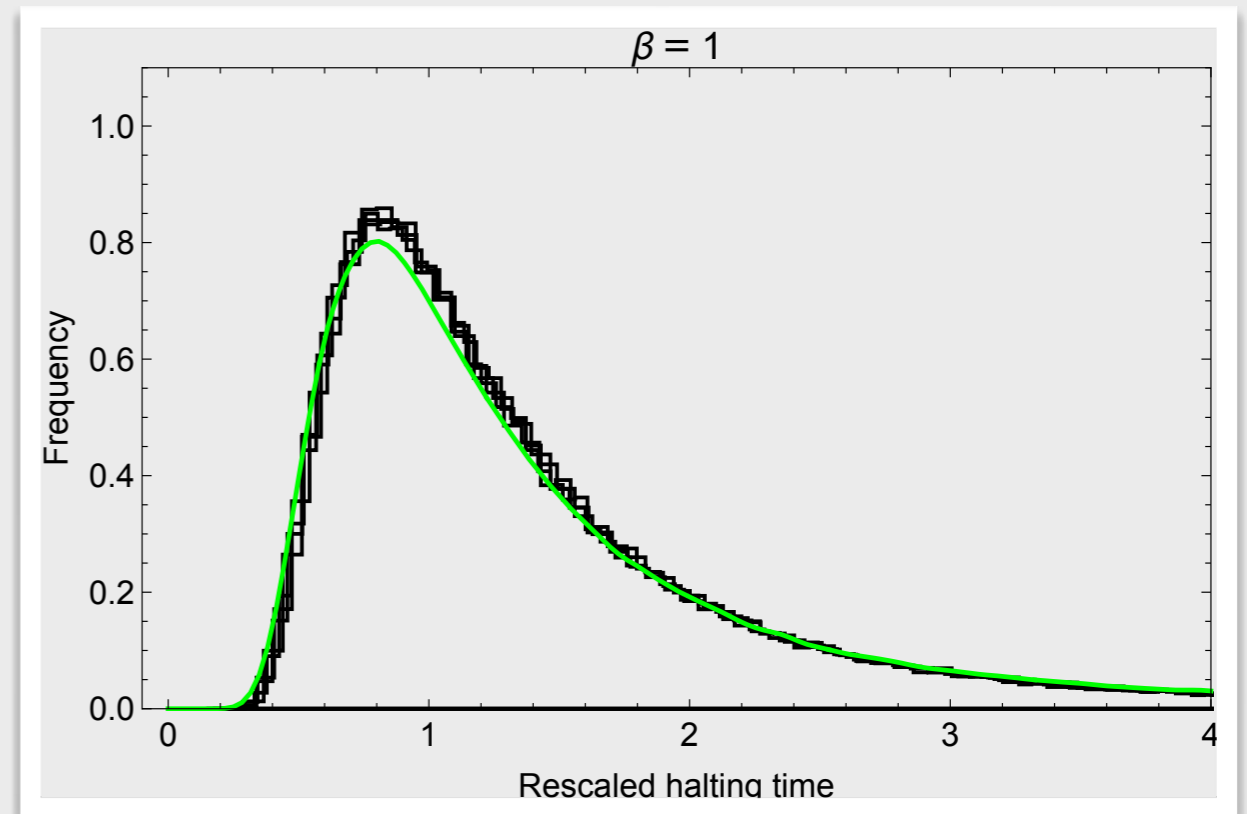
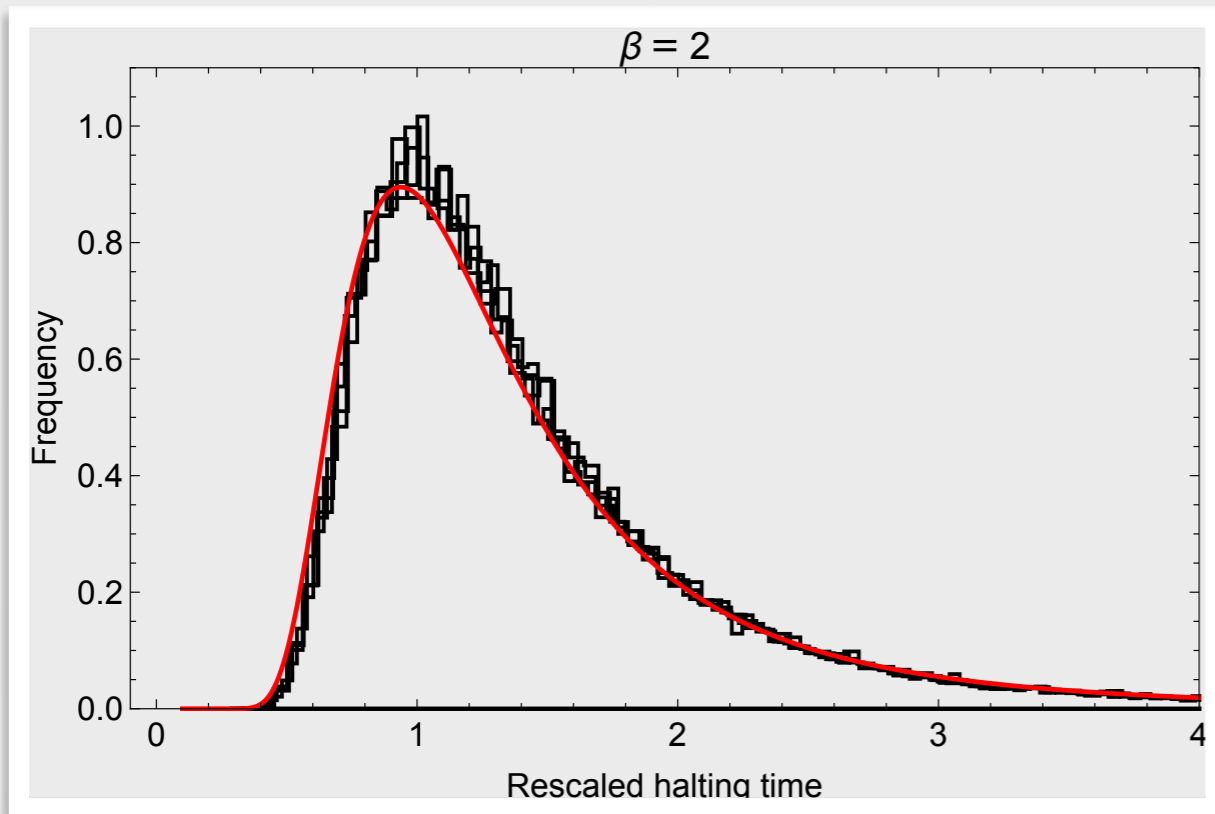


Concentration occurs if $\frac{\epsilon^{-1}}{\sqrt{M}} \ll 1$



A better model of the halting time for the power method?

$N = 500$



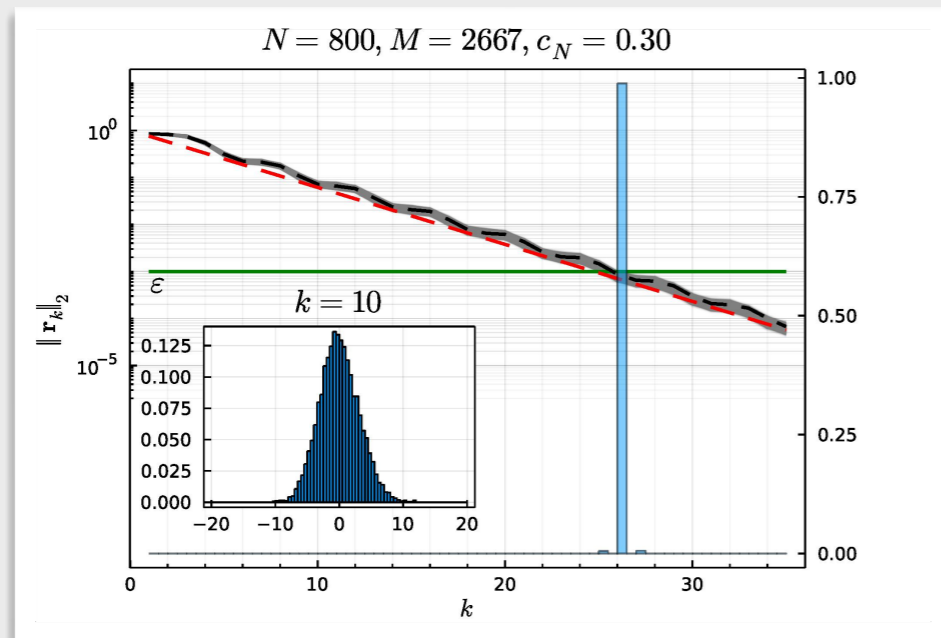
Theorem (Deift and T). Let W be a real ($\beta = 1$) or complex ($\beta = 2$) SCM and let \mathbf{b} be a random unit vector independent of W . Assuming $\epsilon \leq N^{-5/3-\sigma}$, $\sigma > 0$,

$$\lim_{N \rightarrow \infty} \mathbb{P} \left(\frac{T(W, \mathbf{b}, \epsilon)}{2^{-7/6} \lambda_+^{1/3} d^{-1/2} N^{2/3} (\log \epsilon^{-1} - 2/3 \log N)} \leq t \right) = F_{\beta}^{\text{gap}}(t).$$

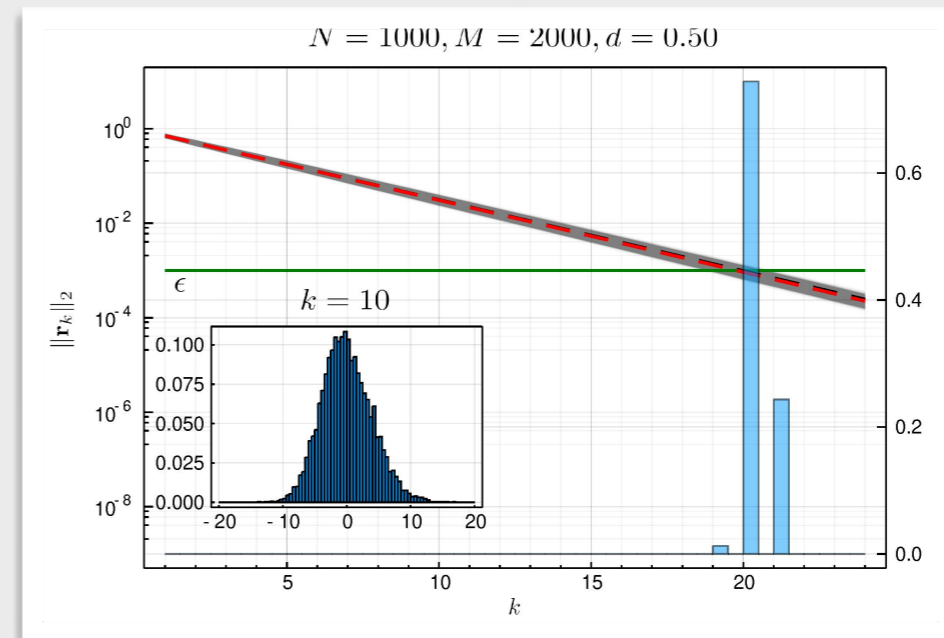
The concentration of the output of algorithms applied to random matrices varies.



Lanczos in finite precision on random matrices



Stabilized CG



Double precision computation

We all know that Lanczos is notoriously unstable.

Once a Ritz value effectively converges, Lanczos typically goes unstable.

But if the Ritz values are behaving like zeros of orthogonal polynomials on a single interval, then the fact that they interlace means they never really converge.

So maybe Lanczos is stable?

C.C. Paige. Accuracy and effectiveness of the Lanczos algorithm for the symmetric eigenproblem. Linear Algebra and its Applications, 34:235–258, dec 1980

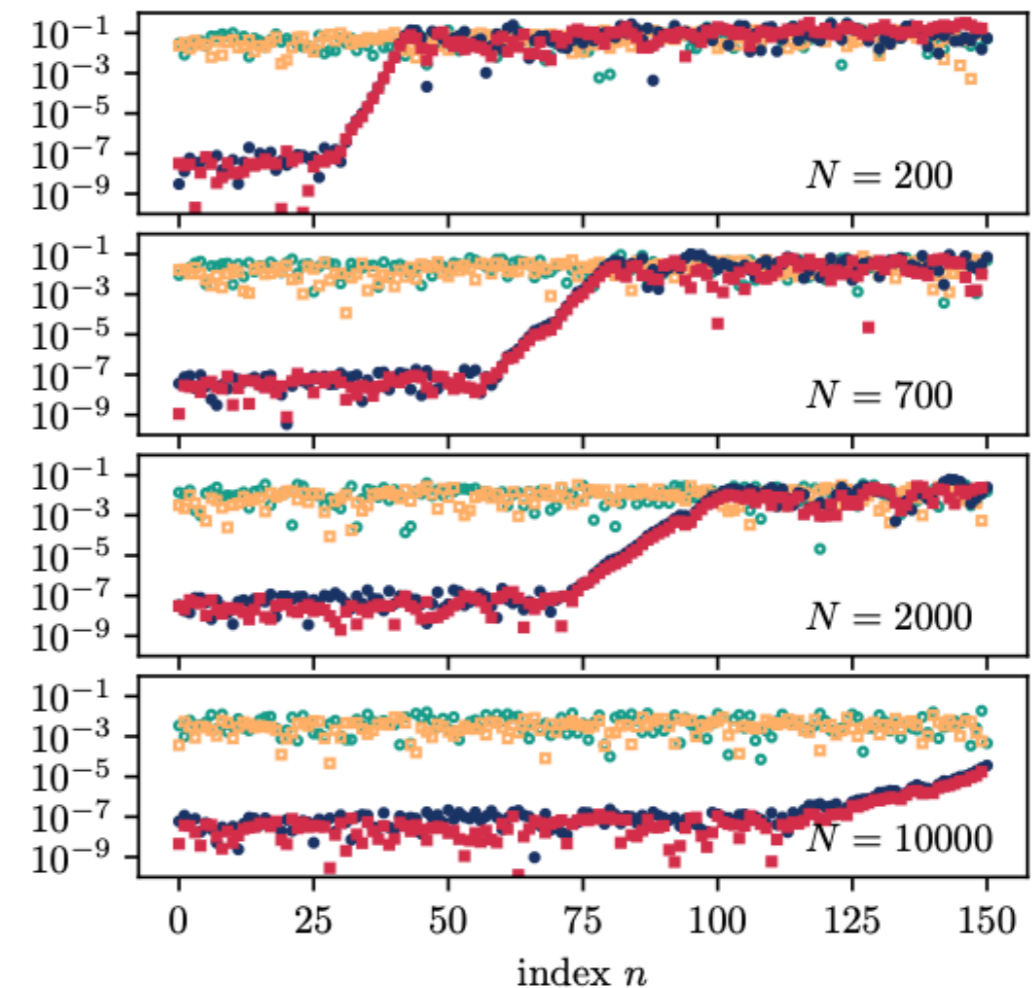


Lanczos is forward/backward stable on classical random matrix ensembles

If $\mu_{W,b} \rightarrow \mu$ and μ is regular in the sense that its orthogonal polynomials are well-behaved, then Lanczos is forward/backward stable.

This applies to the Wishart distribution with overwhelming probability.

Lanczos on a spiked covariance model immediately demonstrates the instability.



(b) Forward error of recurrence coefficients $|\alpha_n - \bar{\alpha}_n|$ (\bullet) and $|\beta_n - \bar{\beta}_n|$ (\blacksquare) and distance to limiting values $|0 - \bar{\alpha}_n|$ (\circ) and $|1/2 - \bar{\beta}_n|$ (\square).

Running algorithms on random matrices may highlight important aspects but hide others!

Motivation to understand more and more random matrix models.



Stability of Lanczos

For an admissible symmetric matrix $A \in \mathbb{R}^{N \times N}$ and a unit vector $\mathbf{b} \in \mathbb{R}^N$ let $T_k(A, \mathbf{b})$ and $\overline{T}_k(A, \mathbf{b})$ be the $k \times k$ Lanczos matrices at iteration k .

There exists $\alpha > 0$ be such that if $\epsilon_{\text{mach}} = O(k^{-\alpha})$ then there exists $\gamma > \beta > 0, \delta > 0$ such that:

1. $T_k(A, \mathbf{b}_*) = \overline{T}_k(A, \mathbf{b})$ where $\|\mathbf{b}_* - \mathbf{b}\| = O(k^{-\gamma})$, and
2. $\|T_k(A, \mathbf{b}) - \overline{T}_k(A, \mathbf{b})\| = O(k^{-\beta} + N^{-\delta})$.



Some details

For CG, we have, approximately,

$$\|\mathbf{e}_k(W, \mathbf{b})\|_W \sim c \left| \int p_k(\lambda; \mu_{W, \mathbf{b}}) \frac{\mu_{W, \mathbf{b}}(d\lambda)}{\lambda} \right| = c \underbrace{\left| \int p_k(\lambda; \mu_{\text{MP}}) \frac{\mu_{\text{MP}}(d\lambda)}{\lambda} \right|}_{\frac{d^{k/2}}{\sqrt{1-d}}} \left(1 + O\left(\sqrt{\frac{k}{M}}\right) \right),$$

where $p_k(x; \mu)$ is $L^2(\mu)$ normalized. It turns out that perturbations of this quantity, preserve its exponential scale because modified moments are well-behaved.

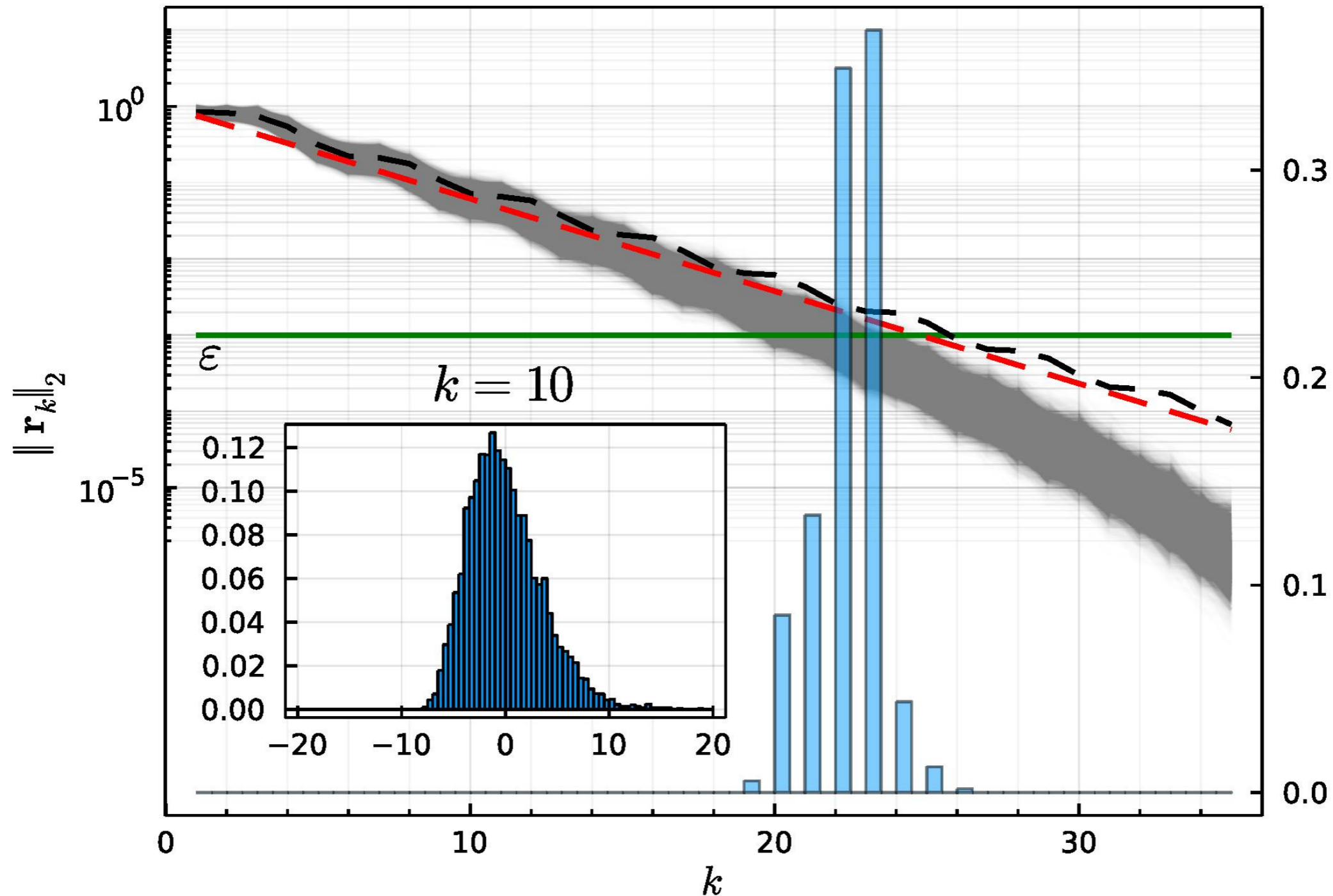
For the power method

$$E_k(W, \mathbf{b}) = \left(\underbrace{(1 + \sqrt{d})^2 - \frac{\int \lambda^{2k-1} \mu_{\text{MP}}(d\lambda)}{\int \lambda^{2k-2} \mu_{\text{MP}}(d\lambda)}}_{O(k^{-1})} + O\left(k^{-2/3} \sqrt{\frac{\log k}{M}}\right) \right).$$



Spiked covariance with multiple bulk components

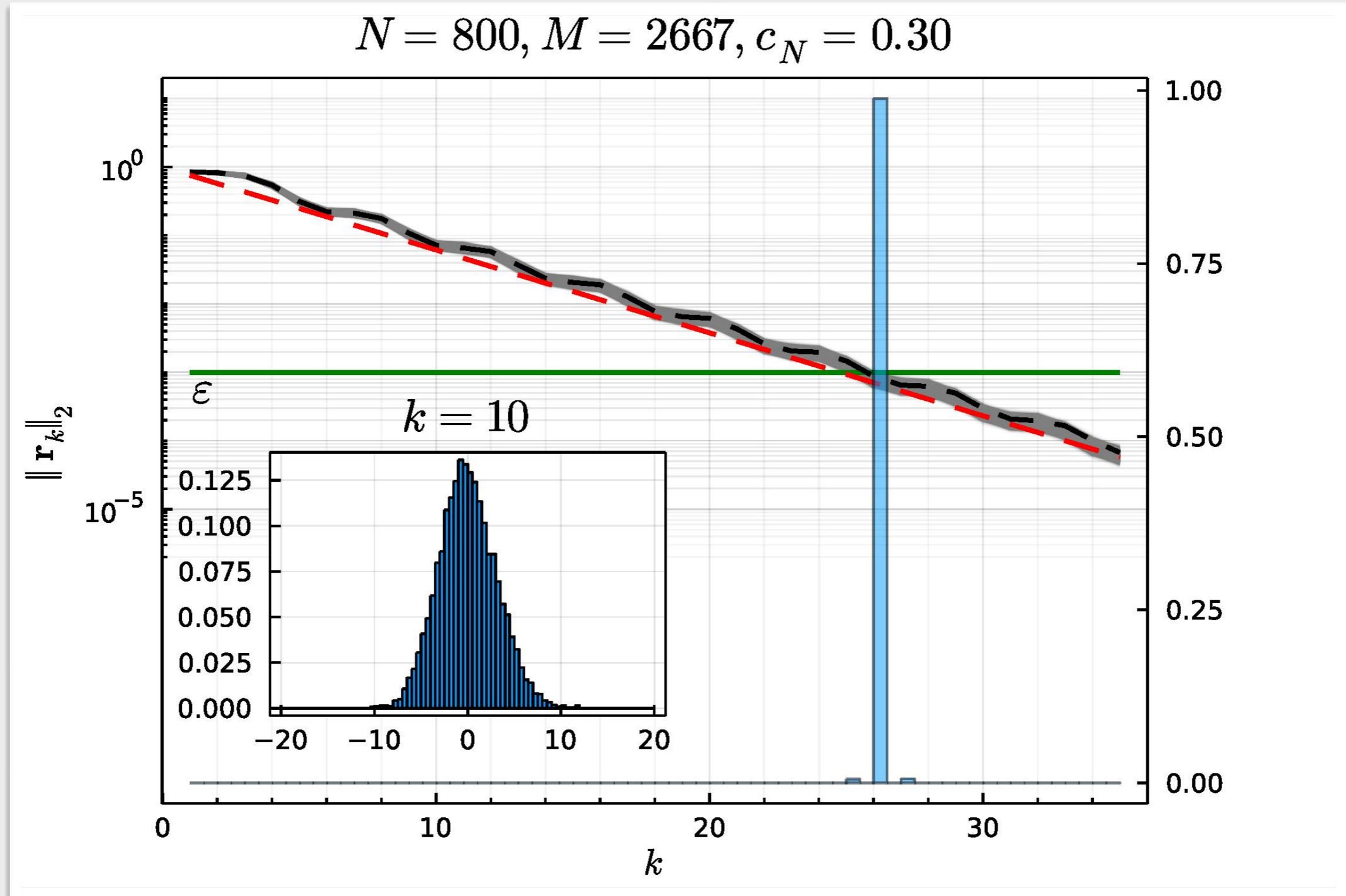
$$N = 50, M = 167, c_N = 0.30$$



X Ding and T T. A Riemann–Hilbert approach to the perturbation theory for orthogonal polynomials: Applications to numerical linear algebra and random matrix theory. [arXiv preprint 2112.12354](#), pages 1–77, 2021



Spiked covariance with multiple bulk components



X Ding and T T. A Riemann–Hilbert approach to the perturbation theory for orthogonal polynomials: Applications to numerical linear algebra and random matrix theory. [arXiv preprint 2112.12354](#), pages 1–77, 2021



The CLT

For general sample covariance matrices, we identify Gaussian processes $\mathcal{G}^{(1)} = (\mathcal{G}_j^{(1)})_{j \geq 1}$, $\mathcal{G}^{(2)} = (\mathcal{G}_j^{(2)})_{j \geq 1}$ so that for $\mathfrak{d} = N/M \rightarrow d \in (0, 1]$

$$\sqrt{\frac{\beta M}{2}} \left(\|\mathbf{e}_k(W, \mathbf{b})\|_W^2 - \frac{\mathfrak{d}^k}{1 - \mathfrak{d}} \right)_{k \geq 1} \xrightarrow{\text{dist.}} \mathcal{G}^{(1)}, \quad d \neq 1,$$

$$\sqrt{\frac{\beta M}{2}} \left(\|\mathbf{r}_k(W, \mathbf{b})\|_2^2 - \mathfrak{d}^k \right)_{k \geq 1} \xrightarrow{\text{dist.}} \mathcal{G}^{(2)},$$

in the sense of convergence of finite-dimensional marginals.

In particular, this implies

$$\sqrt{\frac{\beta M}{2}} \left(\|\mathbf{r}_k(W, \mathbf{b})\|_2^2 - \mathfrak{d}^k \right) \xrightarrow{\text{dist.}} \mathcal{N}(0, \sigma_{k,r}^2), \quad \sigma_{k,r}^2 = k d^{2k} \left(1 + \frac{1}{d} \right).$$

Importantly, we have universality in the sense that this holds for a wide class of sample covariance matrices W (with trivial covariance) that goes far beyond the Gaussian case of Wishart (LUE & LOE).

