# RANDOMIZED BLOCK KRYLOV METHODS
# FOR APPROXIMATING EXTREME EIGENVALUES

JOEL A. TROPP

ABSTRACT. Randomized block Krylov subspace methods form a powerful class of algorithms for computing the extreme eigenvalues of a symmetric matrix or the extreme singular values of a general matrix. The purpose of this paper is to develop new theoretical bounds on the performance of randomized block Krylov subspace methods for these problems. For matrices with polynomial spectral decay, the randomized block Krylov method can obtain an accurate spectral norm estimate using only a *constant* number of steps (that depends on the decay rate and the accuracy). Furthermore, the analysis reveals that the behavior of the algorithm depends in a delicate way on the block size. Numerical evidence confirms these predictions.

## 1. MOTIVATION AND MAIN RESULTS

Randomized block Krylov methods have emerged as a powerful tool for spectral computation and matrix approximation [RST09, HMT11, MRT11, HMST11, MM15, WZZ15, DIKMI18, DI19, YGL18, MT20]. At present, our understanding of these methods is more rudimentary than our understanding of simple Krylov subspace methods in either the deterministic or random setting. Our aim is to develop *a priori* bounds that better explain the remarkable performance of randomized block Krylov methods.

This paper focuses on the most basic questions. How well can we estimate the maximum (or minimum) eigenvalue of a symmetric matrix using a randomized block Krylov method? How well can we estimate the maximum (or minimum) singular value of a general matrix? We present the first detailed analysis that addresses these questions. The argument streamlines and extends the influential work [KW92] of Kucziński & Woźniakowski on the simple Krylov method with one randomized starting vector. We discover that increasing the block size offers some clear advantages over the simple method.

Randomized block Krylov methods have also been promoted as a tool for low-rank matrix approximation. Theoretical analysis of these algorithms appears in several recent papers [MM15, DIKMI18, DI19, YGL18]. In an unpublished companion paper [Tro18], we complement these works with a more detailed analysis that offers several new insights. The mathematical approach is different from the case of extreme eigenvalues, and it involves additional ideas from the analysis of the randomized SVD algorithm [HMT11]. In collaboration with Robert J. Webber, we are currently preparing a review article [TW21] with a definitive treatment of randomized block Krylov methods for low-rank matrix approximation.

1.1. **Block Krylov Methods for Computing the Maximum Eigenvalue.** Let us begin with a mathematical description of a block Krylov method for estimating the maximum eigenvalue of a symmetric matrix. See Section 1.1.6 for a brief discussion about implementations.

1.1.1. *Block Krylov Subspaces.* Suppose that we are given a symmetric matrix $A \in \mathbb{R}^{n \times n}$. Choose a *test matrix* $B \in \mathbb{R}^{n \times \ell}$, where the number $\ell$ is called the *block size*. Select a parameter $q \in \mathbb{N}$ that controls the *depth* of the Krylov subspace. In concept, the block Krylov method constructs the matrix

$$S_q(A; B) := \begin{bmatrix} B & AB & A^2B & \dots & A^qB \end{bmatrix} \in \mathbb{R}^{n \times (q+1)\ell}. \tag{1.1}$$

The range $K_q$ of the matrix $S_q$ is called a *block Krylov subspace*:

$$K_q(A; B) := \text{range}(S_q) \subset \mathbb{R}^n. \tag{1.2}$$

The block Krylov subspace admits an alternative representation in terms of polynomials:

$$K_q(A; B) = \text{span}\left\{\text{range}\left(\varphi(A)B\right) : \varphi \in \mathscr{P}_q\right\}, \tag{1.3}$$

where $\mathscr{P}_q$ is the linear space of real polynomials with degree at most $q$.

**Warning 1.1** (Notation)**.** Be aware that our notation differs slightly from the traditional usage. In $S_q(A; B)$ and in $K_q(A; B)$, the highest matrix power is $A^q$, rather than the usual $A^{q-1}$.

1.1.2. *Invariance Properties of Krylov Subspaces.* Krylov subspaces have remarkable invariance properties that help explain their computational value.

- The block Krylov subspace only depends on the range of the test matrix:

$$K_q(A; BT) = K_q(A; B) \quad \text{for all nonsingular } T \in \mathbb{R}^{\ell \times \ell}.$$

- The block Krylov subspace co-varies with the orientation of the matrices:

$$K_q(UAU^*; UB) = U K_q(A; B) \quad \text{for all orthogonal } U \in \mathbb{R}^{n \times n}.$$

- The block Krylov subspace is invariant under "affine" transformations of the input matrix:

$$K_q(\alpha A + \beta \mathbf{I}; B) = K_q(A; B) \quad \text{for all } \alpha, \beta \in \mathbb{R}.$$

These facts follow directly from the definition (1.1)–(1.2) of the Krylov subspace and the representation (1.3) using polynomials. For example, see [Par98, Sec. 12.2.2] for the simple case $\ell = 1$.

1.1.3. *Computing Maximum Eigenvalues.* Block Krylov subspaces support a wide range of matrix computations. The core idea is to compress the input matrix to the Krylov subspace and to perform calculations on the (small) compressed matrix. In other words, Krylov methods belong to the class of Ritz–Galerkin methods; see [Lan50, Pai71, Saa80, Par98, Saa11].

In particular, we can obtain an estimate $\xi_{\max}(A; B; q)$ for the maximum eigenvalue $\lambda_{\max}(A)$ of the input matrix by maximizing the Rayleigh quotient of $A$ over the block Krylov subspace $K_q(A; B)$:

$$\xi_{\max}(A; B; q) := \max_{v \in K_q(A;B)} \frac{v^* A v}{v^* v}. \tag{1.4}$$

The symbol $^*$ denotes the transpose of a matrix or vector, and we instate the convention that $0/0 = 0$. We may suppress the dependence of $\xi_{\max}$ on $A$, $B$, or $q$ when they are clear from context.

The Rayleigh–Ritz theorem [Bha97, Cor. III.2.1] implies that the maximum eigenvalue estimate (1.4) always satisfies

$$\lambda_{\min}(A) \le \xi_{\max}(A; B; q) \le \lambda_{\max}(A). \tag{1.5}$$

The goal of our analysis is to understand how well $\xi_{\max}(A; B; q)$ *approximates* the maximum eigenvalue $\lambda_{\max}(A)$ as a function of the block size $\ell$ and the depth $q$ of the Krylov space.

Provided that $\xi_{\max} \approx \lambda_{\max}$, any vector $v_{\max}$ that maximizes the Rayleigh quotient in (1.4) has a component in the invariant subspace of $A$ spanned by the large eigenvalues. More precisely,

$$\left(\frac{\|P_{>\lambda} v_{\max}\|}{\|v_{\max}\|}\right)^2 \ge \frac{\xi_{\max} - \lambda}{\lambda_{\max} - \lambda} \quad \text{for all } \lambda < \xi_{\max}.$$

We have written $\|\cdot\|$ for the $\ell_2$ norm and $P_{>\lambda}$ for the spectral projector onto the invariant subspace of $A$ associated with the eigenvalues that strictly exceed $\lambda$.

1.1.4. *Invariance Properties of the Eigenvalue Estimate.* The eigenvalue estimate $\xi_{\max}(A; B; q)$ inherits some invariance properties from the block Krylov subspace. These facts can help us develop effective implementations of the algorithm and to analyze their performance.

- For fixed depth $q$, the estimate only depends on the range of the test matrix $B$:

$$\xi_{\max}(A; BT) = \xi_{\max}(A; B) \quad \text{for all nonsingular } T \in \mathbb{R}^{\ell \times \ell}. \tag{1.6}$$

- For fixed depth $q$, the estimate does not depend on the orientation of $A$ and $B$ in the sense that

$$\xi_{\max}(UAU^*; UB) = \xi_{\max}(A; B) \quad \text{for all orthogonal } U \in \mathbb{R}^{n \times n}. \tag{1.7}$$

- For fixed depth $q$, the estimate covaries with *increasing* "affine" transformations of $A$:

$$\xi_{\max}(\alpha A + \beta \mathbf{I}; B) = \alpha \xi_{\max}(A; B) + \beta \quad \text{for all } \alpha \geq 0 \text{ and } \beta \in \mathbb{R}. \tag{1.8}$$

These results all follow immediately from the invariance properties of the Krylov subspace (Section 1.1.2) and the definition (1.4) of the eigenvalue estimate. See [KW92] for the simple case $\ell = 1$.

1.1.5. *A Random Test Matrix.* To ensure that we can estimate the maximum eigenvalue of an *arbitrary* symmetric input matrix $A$, we draw the test matrix $B$ that generates the Krylov subspace *at random*.

How should we select the distribution? Observe that the eigenvalue estimate $\xi_{\max}(A; B; q)$ only depends on the range of the test matrix $B$ because of the property (1.6). Furthermore, the property (1.7) shows that the eigenvalue estimate is invariant under rotations. Therefore, we can choose any random test matrix whose range has a rotationally invariant distribution. This idea is extracted from [KW92]; see Section 4.1 for the justification.

We will consider a standard normal test matrix $\Omega \in \mathbb{R}^{n \times \ell}$. That is, the entries of $\Omega$ are statistically independent Gaussian random variables, each with mean zero and variance one. It is well known that the range of this random matrix has a rotationally invariant distribution. The goal of this paper is to study the behavior of the random eigenvalue estimate $\xi_{\max}(A; \Omega; q)$.

**Remark 1.2** (Other Test Matrices)**.** Our analysis and detailed results depend heavily on the choice of a standard normal test matrix $\Omega$. In practice, we can achieve similar empirical behavior from test matrices with "structured" distributions that require less storage or that have fast matrix–vector multiplies. For randomized Krylov subspace methods, the computational benefit of using a structured test matrix is limited because we need to perform repeated multiplications with the input matrix $A$ to generate the Krylov matrix; cf. (1.1). See [HMT11, Secs. 4.6 and 7.4] or [MT20, Secs. 9 and 11.5] for some discussion about other random test matrices.

1.1.6. *Implementation.* For completeness, we describe the simplest stable implementation of a block Krylov method for computing the largest eigenvalue of a symmetric matrix. See Algorithm 1 for pseudo-code; this approach is adapted from [RST09, HMST11]. Here is a summary of the costs:

- A total of $q$ matrix–matrix multiplies between $A$ and an $n \times \ell$ matrix, plus another multiplication between $A$ and an $n \times (q+1)\ell$ matrix. The arithmetic cost depends on whether the matrix $A$ supports a fast multiplication operation. For example, the algorithm is more efficient when $A$ is sparse.
- Repeated orthogonalization of $(q+1)\ell$ vectors of length $n$ with total cost $\mathcal{O}(q^2 \ell^2 n)$ operations.
- Solution of a (block-tridiagonal) maximum eigenvalue problem at a cost of $\mathcal{O}(q^2 \ell^2)$ operations.
- Maintenance of the matrix $S_q$, which requires $\mathcal{O}(q \ell n)$ units of storage.

More refined algorithms can reduce the resource usage markedly; see [MT20, Sec. 11.7] or [TW21] for discussion.

As noted, Krylov subspace methods are particularly valuable when we have an efficient procedure for computing matrix–vector products with $A$. On many contemporary computer architectures, the cost of performing a product with several vectors is similar to the cost of a product with a single vector. In this setting, we can take advantage of the improved behavior of a block method almost for free.

---

**Algorithm 1** A block Krylov method for computing the largest eigenvalue of a symmetric matrix

---

**Input:** Symmetric $n \times n$ matrix $A$; block size $\ell$; depth $q$
**Output:** Estimate $\xi$ of largest eigenvalue

---

```
 1  function BLOCKKRYLOVMAXEIG(A, ℓ, q)
 2     [B₀, ~] ← qr_econ(randn(size(A,1),ℓ))          ▷ Draw n × ℓ random orthonormal matrix
 3     S₀ = B₀                                         ▷ Orthonormal basis for Krylov subspace
 4     for t ← 1,2,3,…,q do
 5        Bₜ ← ABₜ₋₁                                   ▷ Form next block of Krylov matrix by multiplication
 6        Bₜ ← Bₜ − Sₜ₋₁(S*ₜ₋₁Bₜ)                      ▷ Orthogonalize against subspace by double GS
 7        Bₜ ← Bₜ − Sₜ₋₁(S*ₜ₋₁Bₜ)                      ▷ Sic!
 8        [Bₜ, ~] ← qr_econ(Bₜ)                        ▷ Orthogonalize what's left
 9        Sₜ ← [Sₜ₋₁, Bₜ]                              ▷ Extend basis with new block
10     H ← S*q(ASq)                                    ▷ Form block-tridiagonal Rayleigh quotient matrix
11     [ξ, v] ← maxeig(H)                              ▷ Find maximum eigenpair
```

---

We refer to the books [Par98, BDD$^+$00, Saa11, GVL13] and the paper [HMST11] for more discussion and references.

1.2. **The Role of the Spectrum.** Owing to invariance, the theoretical behavior of the eigenvalue estimate $\xi_{\max}(A; \Omega; q)$ depends only on the spectrum of the input matrix $A$. In this section, we develop this idea further and introduce some spectral features that affect the performance of the eigenvalue estimate.

**Warning 1.3** (Numerical Behavior)**.** The numerical performance of a (block) Krylov method can be very complicated; for example, see [Par98, Chap. 13] or [Meu06]. The current paper does not address numerics.

1.2.1. *Invariance Properties of the Random Eigenvalue Estimate.* The random estimate $\xi_{\max}(A; \Omega; q)$ of the maximum eigenvalue has several invariance properties that allow us to simplify the analysis.

First, the rotation invariance (1.7) of the eigenvalue estimate and the rotational invariance of the range of $\Omega$ imply that

$$\xi_{\max}(A; \Omega) \sim \xi_{\max}(\Lambda; \Omega) \quad \text{where } A = U\Lambda U^* \text{ is an eigenvalue factorization.} \tag{1.9}$$

The symbol $\sim$ signifies equality of distribution for two random variables. In other words, the maximum eigenvalue estimate depends only on the eigenvalues of the input matrix—but not the eigenvectors.

Second, owing to the affine covariance property (1.8), the eigenvalue estimate $\xi_{\max}(A; \Omega; q)$ only depends on the "shape" of the spectrum of $A$, but not its location or scale. As a consequence, we must assess the behavior of the eigenvalue estimate in terms of spectral features that are affine invariant.

1.2.2. *Spectral Features of the Input Matrix.* To express the results of our analysis, we introduce terminology for some spectral features of the symmetric matrix $A \in \mathbb{R}^{n \times n}$. First, let us instate compact notation for the eigenvalues of $A$:

$$a_i := \lambda_i(A) \quad \text{for } i = 1,\ldots,n, \quad \text{and} \quad a_{\max} := a_1 \geq a_2 \geq \cdots \geq a_n =: a_{\min}.$$

The map $\lambda_i(\cdot)$ returns the $i$th largest eigenvalue of a symmetric matrix.

Next, we define some functions of the eigenvalue spectrum:

- The *spectral range* $\rho$ of the input matrix is the distance between the extreme eigenvalues. That is, $\rho := a_{\max} - a_{\min}$.
- The *spectral gap* $\gamma$ is the relative difference between the maximum eigenvalue and the next distinct eigenvalue:

$$\gamma := \frac{a_{\max} - a_{m+1}}{a_{\max} - a_{\min}} \quad \text{where} \quad a_{\max} = a_m > a_{m+1}. \tag{1.10}$$

If $A$ is a multiple of the identity, then $\gamma = 0$. Note that $\gamma \in [0,1]$.

- Let $\nu$ be a nonnegative number. The $\nu$-*stable rank* is a continuous measure of the "dimension" of the range of $A - a_{\min}\mathbf{I}$ that reflects how quickly the spectrum decays. It is defined as

$$\mathrm{srk}(\nu) := \sum_{i=1}^{n} \left( \frac{a_i - a_{\min}}{a_{\max} - a_{\min}} \right)^{2\nu}. \tag{1.11}$$

  If $A$ is a multiple of the identity, we define $\mathrm{srk}(\nu) = 0$. Otherwise, $1 \le \mathrm{srk}(\nu) \le \mathrm{rank}(A - a_{\min}\mathbf{I}) \le n - 1$. When the eigenvalues of $A$ decay at a polynomial rate, the stable rank can be much smaller than the rank for an appropriate choice of $\nu$. See Section 2.1 for some illustrations.
- Let $\xi$ be any estimate for the largest eigenvalue $a_{\max}$ of the input matrix $A$. We measure the error in the estimate relative to the spectral range:

$$\mathrm{err}(\xi) := \frac{a_{\max} - \xi}{a_{\max} - a_{\min}}. \tag{1.12}$$

  The relative error in the Krylov estimate $\xi = \xi_{\max}(A; B; q)$ falls in the interval $[0, 1]$ because of (1.5).

The spectral gap, the stable rank, and the error measure are all invariant under increasing affine transformations of the spectrum of $A$. We suppress the dependence of these quantities on the input matrix $A$, unless emphasis is required.

**Remark 1.4** (History)**.**  Concepts related to the $\nu$-stable rank originally appeared in the analysis literature [BT87, BT91], and they now play an important role in randomized linear algebra [MT20, Sec. 2.6].

1.2.3. *An Example.*  Let us illustrate these concepts with a canonical example. A common task in computational physics is to compute the *minimum* eigenvalue of (the discretization of) an elliptic operator.

Consider the standard $n$-point second-order finite-difference discretization $L \in \mathbb{R}^{n \times n}$ of the one-dimensional Laplacian on $[0, 1]$ with homogeneous boundary conditions. Writing $h = 1/(n + 1)$, the eigenvalues of the matrix $L$ are

$$\lambda_j(L) = \frac{2}{h^2} \left[ 1 + \cos(\pi j h) \right] \quad \text{for } j = 1, \dots, n.$$

The spectral gap between the smallest and second smallest eigenvalue is proportional to $h^2$. The spectrum has essentially no decay, which is visible in the fact that $\mathrm{srk}(-L; 1)$ is proportional to $1/h$. The tiny spectral gap and large stable rank suggest that the maximum eigenvalue problem for $-L$ will be challenging.

A natural remedy is to attempt to compute the maximum eigenvalue of the inverse $L^{-1}$. Independent of $n$, the spectral gap between the largest and second largest eigenvalue satisfies $\gamma(L^{-1}) \approx 0.75$. The stable rank $\mathrm{srk}(L^{-1}; 1) \approx 1.1$. The large spectral gap and small stable rank suggest that the maximum eigenvalue problem for $L^{-1}$ will be quite easy, regardless of the dimension.

1.3. **Matrices with Few Distinct Eigenvalues.**  Before continuing, we must address an important special case. When the input matrix has few distinct eigenvalues, the block Krylov method computes the maximum eigenvalue of the matrix perfectly.

**Proposition 1.5** (Randomized Block Krylov: Matrices with Few Eigenvalues)**.**  *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric matrix. Fix the block size $\ell \ge 1$ and the depth $q \ge 0$ of the block Krylov subspace. Draw a standard normal matrix $\Omega \in \mathbb{R}^{n \times \ell}$. If $A$ has $q + 1$ or fewer distinct eigenvalues, then $\mathrm{err}(\xi_{\max}(A; \Omega; q)) = 0$ with probability one.*

This type of result is well known (e.g., see [KW92]), but we include a short proof in Section 5.

1.4. **Matrices without a Spectral Gap.**  Our first result gives probabilistic bounds for the maximum eigenvalue estimate $\xi_{\max}(A; \Omega; q)$ without any additional assumptions. In particular, it does not require a lower bound on the spectral gap $\gamma$.

**Theorem 1.6** (Randomized Block Krylov: Maximum Eigenvalue Estimate)**.**  *Instate the following hypotheses.*

- *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric input matrix.*
- *Draw a standard normal test matrix $\Omega \in \mathbb{R}^{n \times \ell}$ with block size $\ell$.*
- *Fix the depth parameter $q \geq 0$, and let $q = q_1 + q_2$ be an arbitrary nonnegative integer partition.*

*We have the following probability bounds for the estimate $\xi_{\max}(A; \Omega; q)$, defined in (1.4), of the maximum eigenvalue of the input matrix.*

*(1) The relative error (1.12) in the eigenvalue estimate satisfies the probability bound*

$$\mathbb{P}\left\{\mathrm{err}(\xi_{\max}(A; \Omega; q)) \geq \varepsilon\right\} \leq 1 \wedge \sqrt{2}\left[8\,\mathrm{srk}(q_1) \cdot \mathrm{e}^{-2(2q_2+1)\sqrt{\varepsilon}}\right]^{\ell/2} \quad \text{for } \varepsilon \in [0, 1]. \tag{1.13}$$

*(2) The expectation of the relative error satisfies*

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(A; \Omega; q)) \leq 1 \wedge \left[\frac{2.70\ell^{-1} + \log(8\,\mathrm{srk}(q_1))}{2(2q_2 + 1)}\right]^2. \tag{1.14}$$

*The symbol $\wedge$ denotes the minimum, and the stable rank, $\mathrm{srk}(\cdot)$, is defined in (1.11). All logarithms are natural.*

The proof of (1.13) appears in Section 6, and the proof of (1.14) appears in Section 7. The experiments in Section 2 support the analysis.

Let us take a moment to explain the content of this result. We begin with a discussion about the role of the second depth parameter $q_2$, and then we explain the role of the first depth parameter $q_1$. We emphasize that the user does not choose a partition $q = q_1 + q_2$; the bounds are valid for all partitions.

For now, fix $q_1$. The key message of Theorem 1.6 is that the relative error satisfies

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(A; \Omega; q_1 + q_2)) \leq \varepsilon$$

once the depth parameter $q_2$ exceeds

$$q_2(\varepsilon) := -\frac{1}{2} + \frac{2.70\ell^{-1} + \log(8\,\mathrm{srk}(q_1))}{4\sqrt{\varepsilon}}.$$

Once the depth $q_2$ attains this level, the probability of error drops off exponentially fast:

$$q_2 = q_2(\varepsilon) + k\varepsilon^{-1/2} \quad \text{implies} \quad \mathbb{P}\left\{\mathrm{err}(\xi_{\max}(A; \Omega; q_1 + q_2)) \geq \varepsilon\right\} \leq \mathrm{e}^{-2k\ell}.$$

In fact, we need $q_2 \geq q_2(1)$ just to ensure that the probability bound is nontrivial.

The most important aspect of this result is that the depth $q_2(\varepsilon)$ scales with $\varepsilon^{-1/2}$, so it is possible to achieve moderate relative error using a block Krylov space with limited depth. In contrast, the randomized power method [KW92] requires the depth $q$ to be proportional to $\varepsilon^{-1}$ to achieve a relative error of $\varepsilon$.

The second thing to notice is that the depth $q_2(\varepsilon)$ scales with $\log(\mathrm{srk}(q_1))$. The stable rank is never larger than the ambient dimension $n$, but it can be significantly smaller—even constant—when the spectrum of the matrix has polynomial decay.

Here is another way to look at these facts. As we increase the depth parameter $q$, the block Krylov method exhibits a burn-in period whose length $q_1 + q_2(1)$ depends on $\mathrm{srk}(q_1)$. While the depth $q_2 \leq q_2(1)$, the algorithm does not make much progress in estimating the maximum eigenvalue. Once the depth satisfies $q_2 \geq q_2(1)$, the expected relative error decreases in proportion to $q_2^{-2}$. In contrast, the power method [KW92] reduces the expected relative error in proportion to $q_2^{-1}$.

We can now appreciate the role of the first depth parameter $q_1$. When the spectrum of the input matrix exhibits polynomial decay, $\mathrm{srk}(q_1)$ is *constant* for an appropriate value of $q_1$ that depends on the decay rate. In this case, the analysis shows that the total burn-in period $q_1 + q_2(1)$ is just $\mathcal{O}(1)$ steps. For example, when the $j$th eigenvalue decays like $1/j$, this situation occurs.

The block size $\ell$ may not play a significant role in determining the average error. But changing the block size has a large effect on the probability of failure (i.e., the event that the relative error exceeds $\varepsilon$). For example, suppose that we increase the block size $\ell$ from one to three. For each increment of $\varepsilon^{-1/2}$ in the depth $q_2$, the failure probability with block size $\ell = 3$ is a factor of $54\times$ smaller than the failure probability with block size $\ell = 1$!

**Remark 1.7** (Prior Work)**.** The simple case $\ell = 1$ in Theorem 1.6 has been studied in the paper [KW92]. Our work introduces two major innovations. First, we obtain bounds in terms of the stable rank, which allows us to mitigate the dimensional dependence that appears in [KW92]. Second, we have obtained precise results for larger block sizes $\ell$, which indicate potential benefits of using block Krylov methods. Our proof strategy is motivated by the work in [KW92], but we have been able to streamline and extend their arguments by using a more transparent random model for the test matrix.

1.5. **Matrices with a Spectral Gap.** Our second result gives probabilistic bounds for the maximum eigenvalue estimate $\xi_{\max}(A; \Omega; q)$ when we have a lower bound for the spectral gap $\gamma$ of the input matrix.

**Theorem 1.8** (Randomized Block Krylov: Maximum Eigenvalue Estimate with Spectral Gap)**.** *Instate the following hypotheses.*

- *Let $A \in \mathbb{R}^{n \times n}$ be a symmetric input matrix with spectral gap $\gamma$, defined in (1.10).*
- *Draw a standard normal test matrix $\Omega \in \mathbb{R}^{n \times \ell}$ with block size $\ell$.*
- *Fix the depth parameter $q \geq 0$, and let $q = q_1 + q_2$ be an arbitrary nonnegative integer partition.*

*We have the following probability bounds for the estimate $\xi_{\max}(A; \Omega; q)$, defined in (1.4), of the maximum eigenvalue of the input matrix.*

*(1) The relative error (1.12) in the eigenvalue estimate satisfies the probability bound*

$$\mathbb{P}\left\{\operatorname{err}(\xi_{\max}(A; \Omega; q)) \geq \varepsilon\right\} \leq 1 \wedge \sqrt{2}\left[\frac{8\operatorname{srk}(q_1)}{\varepsilon} \cdot \mathrm{e}^{-4q_2\sqrt{\gamma}}\right]^{\ell/2} \quad \textit{for } \varepsilon \in (0, 1]. \tag{1.15}$$

*(2) Abbreviate $F := 4\operatorname{srk}(q_1)\mathrm{e}^{-4q_2\sqrt{\gamma}}$. The expectation of the relative error satisfies*

$$\mathbb{E}\operatorname{err}(\xi_{\max}(A; \Omega)) \leq \frac{F}{(\ell - 2) + F} \qquad\qquad (\ell \geq 3); \tag{1.16}$$

$$\mathbb{E}\operatorname{err}(\xi_{\max}(A; \Omega)) \leq (F/2)\log(1 + 2/F) \qquad (\ell = 2); \tag{1.17}$$

$$\mathbb{E}\operatorname{err}(\xi_{\max}(A; \Omega)) \leq 1 \wedge \sqrt{2\pi F} \qquad\qquad (\ell = 1). \tag{1.18}$$

*The symbol $\wedge$ denotes the minimum, and the stable rank is defined in (1.11).*

The proof of (1.15) appears in Sections 6; the proof of (1.16), (1.17), and (1.18) appears in Section 8. The experiments in Section 2 bear out these predictions.

Let us give a verbal summary of what this result means. First of all, we anticipate that Theorem 1.8 will give better bounds than Theorem 1.6 when the spectral gap $\gamma$ exceeds the target level $\varepsilon$ for the error. But both results are valid for all choices of $\gamma$ and $\varepsilon$.

Now, fix the first depth parameter $q_1$. One implication of the spectral gap result is that the relative error satisfies

$$\mathbb{P}\left\{\operatorname{err}(\xi_{\max}(A; \Omega; q)) \leq \varepsilon\right\} \approx 1$$

when the second depth parameter $q_2$ exceeds

$$q_2(\varepsilon; \gamma) := \frac{0.70\ell^{-1} + \log(8\varepsilon^{-1}\operatorname{srk}(q_1))}{4\sqrt{\gamma}}.$$

In this case, the depth $q_2$ only scales with $\log(1/\varepsilon)$, so the block Krylov method can achieve very small relative error for a matrix with a spectral gap. As before, the depth $q_2$ also scales with $\log(\operatorname{srk}(q_1))$, so the dimensional dependence is weak—or even nonexistent if the spectrum has polynomial decay and $q_1$ is sufficiently large.

When the depth $q_2 \geq q_2(\varepsilon; \gamma)$, the error probability drops off quickly:

$$q_2 = q_2(\varepsilon; \gamma) + k\gamma^{-1/2} \quad \text{implies} \quad \mathbb{P}\{\operatorname{err}(\xi_{\max}(A; \Omega)) \geq \varepsilon\} \leq \mathrm{e}^{-2k\ell}.$$

This bound indicates that $\gamma^{-1/2}$ is the scale on which the depth $q_2$ needs to increase to reduce the failure probability by a constant multiple (which depends on the block size).

We discover a new phenomenon when we examine the expectation of the error. On average, to achieve a relative error of $\varepsilon$, it is sufficient that the depth $q_2 \geq q_2'(\varepsilon; \gamma)$, where

$$\ell \geq 3: \quad q_2'(\varepsilon; \gamma) := \frac{\log(4\,\mathrm{srk}(q_1)) + \log(1/\varepsilon) - \log(\ell - 2)}{4\sqrt{\gamma}};$$

$$\ell = 2: \quad q_2'(\varepsilon; \gamma) := \frac{\log(4\,\mathrm{srk}(q_1)) + \log(1/\varepsilon) + \log\log(1/\varepsilon)}{4\sqrt{\gamma}} \quad \text{for all } \varepsilon \leq \text{const};$$

$$\ell = 1: \quad q_2'(\varepsilon; \gamma) := \frac{\log(4\,\mathrm{srk}(q_1)) + 2\log(1/\varepsilon) + \log(2\pi)}{4\sqrt{\gamma}}.$$

In other words, the depth $q_2$ of the block Krylov space needs to be about $\log(4\,\mathrm{srk}(q_1))/(4\sqrt{\gamma})$ before we obtain an average relative error less than one; we can reduce this requirement slightly by increasing the block size $\ell$. But once the depth $q_2$ reaches this level, Theorem 1.8 suggests that the block Krylov method with $\ell \geq 2$ reduces the average error *twice as fast* as the block Krylov method with $\ell = 1$.

**Remark 1.9** (Prior Work)**.** The simple case $\ell = 1$ in Theorem 1.8 has been studied in the paper [KW92]. See Remark 1.7 for a discussion of the improvements we have achieved.

### 1.6. Estimating Minimum Eigenvalues.

We can also use Krylov subspace methods to obtain an estimate $\xi_{\min}(A; B; q)$ for the *minimum* eigenvalue of a symmetric matrix $A$. Conceptually, the simplest way to accomplish this task is to apply the Krylov subspace method to the negation $-A$. The minimum eigenvalue estimate takes the form

$$\xi_{\min}(A; B; q) := -\xi_{\max}(-A; B; q).$$

Owing to (1.5), this estimate is never smaller than $\lambda_{\min}(A)$.

It is straightforward to adapt Theorems 1.6 to obtain bounds for the minimum eigenvalue estimate with a random test matrix $\mathbf{\Omega}$. In particular, we always have the bound

$$\mathbb{E}\left[\frac{\xi_{\min}(A; \mathbf{\Omega}; q) - a_{\min}}{a_{\max} - a_{\min}}\right] \leq \left[\frac{2.70\ell^{-1} + \log(8\,\mathrm{srk}(-A; q_1))}{2(2q_2 + 1)}\right]^2.$$

In this context, the stable rank takes the form

$$\mathrm{srk}(-A; \nu) = \sum_{i=1}^{n} \left(\frac{a_{\max} - a_i}{a_{\max} - a_{\min}}\right)^{2\nu}.$$

See Section 1.4 for discussion of this type of bound.

We can also use Theorem 1.8 to obtain results in terms of the spectral gap. The spectral gap $\gamma(-A)$ is the magnitude of the difference between the smallest two eigenvalues of $A$, relative to the spectral range. For example, when the block size $\ell = 3$, we have the bound

$$\ell = 3: \qquad \mathbb{E}\left[\frac{\xi_{\min}(A; \mathbf{\Omega}; q) - a_{\min}}{a_{\max} - a_{\min}}\right] \leq 4\,\mathrm{srk}(-A; q_1) \cdot \mathrm{e}^{-4q_2\sqrt{\gamma(-A)}}.$$

See Section 1.5 for discussion of this type of bound.

**Remark 1.10** (Inverse Iterations)**.** Suppose that we have a routine for applying the *inverse* of a positive-semidefinite matrix $A$ to a vector. In this case, we can apply the Krylov method to the inverted matrix $A^{-1}$ to estimate the minimum eigenvalue. This approach is often more powerful than applying the Krylov method to $-A$. For example, inverse iteration is commonly used for discretized elliptic operators; the discussion in Section 1.2.3 supports this approach.

### 1.7. Estimating Singular Values.

We now arrive at the problem of estimating the spectral norm of a general matrix $C \in \mathbb{R}^{n \times m}$ using Krylov subspace methods. Assuming that $n \leq m$, we can apply the block Krylov method to the square $CC^*$. This yields an estimate $\xi_{\max}(CC^*; B; q)$ for the *square* of the spectral norm of $C$. If $n > m$, we can just as well work with the other square $C^*C$.

Theorems 1.6 and 1.8 immediately yield error bounds for the random test matrix $\boldsymbol{\Omega}$. In particular, we always have the bound

$$\mathbb{E}\left[\frac{\|\boldsymbol{C}\|^2 - \xi_{\max}(\boldsymbol{CC}^*; \boldsymbol{\Omega}; q)}{\|\boldsymbol{C}\|^2}\right] \leq \left[\frac{2.70\ell^{-1} + \log(8\operatorname{srk}(\boldsymbol{CC}^*; q_1))}{2(2q_2 + 1)}\right]^2.$$

We also have a bound in terms of the spectral gap $\gamma(\boldsymbol{CC}^*)$, which is the difference between the squares of the largest two distinct singular values, relative to the spectral range. For block size $\ell = 3$, we have

$$\mathbb{E}\left[\frac{\|\boldsymbol{C}\|^2 - \xi_{\max}(\boldsymbol{CC}^*; \boldsymbol{\Omega})}{\|\boldsymbol{C}\|^2}\right] \leq 4\operatorname{srk}(\boldsymbol{CC}^*; q_1) \cdot \mathrm{e}^{-4q_2\sqrt{\gamma(\boldsymbol{CC}^*)}}.$$

In this case, it is natural to bound the stable rank as

$$\operatorname{srk}(\boldsymbol{CC}^*; 0) = m \wedge n \quad \text{and} \quad \operatorname{srk}(\boldsymbol{CC}^*; \nu) \leq \left(\frac{\|\boldsymbol{C}\|_{4\nu}}{\|\boldsymbol{C}\|}\right)^{4\nu} \quad \text{for } \nu \geq 1.$$

We have written $\|\cdot\|$ for the spectral norm and $\|\cdot\|_p$ for the Schatten $p$-norm for $p \geq 1$.

**Remark 1.11** (Other Approaches)**.** It is also possible to work with "odd" Krylov subspaces $K_q(\boldsymbol{CC}^*; \boldsymbol{CB})$ or $K_q(\boldsymbol{C}^*\boldsymbol{C}; \boldsymbol{C}^*\boldsymbol{B})$, but the analysis requires some modifications.

**Remark 1.12** (Minimum Singular Values)**.** The quantity $\xi_{\min}(\boldsymbol{CC}^*; \boldsymbol{B})$ gives an estimate for the $m$th largest squared singular value of $\boldsymbol{C}$. When $m \leq n$, this is the smallest singular value, which may be zero. It is straightforward to derive results for the estimate using the principles outlined above. We omit the details.

1.8. **Extensions.** With minor modifications, the analysis in this paper can be extended to cover some related situations. First, when the maximum eigenvalue has multiplicity greater than one, the (block) Krylov method converges more quickly. Second, we can extend the algorithm and the analysis to the problem of estimating the largest eigenvalue of an Hermitian matrix (using a complex standard normal test matrix). Third, we can analyze the behavior of the randomized (block) power method. For brevity, we omit the details.

## 2. NUMERICAL EXPERIMENTS

Krylov subspace methods exhibit complicated behavior because they are implicitly solving a polynomial optimization problem. Therefore, we do not expect a reductive theoretical analysis to capture all the nuances of their behavior. Nevertheless, by carefully choosing input matrices, we can witness the phenomena that the theoretical analysis predicts.

2.1. **Experimental Setup.** We implemented the randomized block Krylov method in MATLAB 2019a. The code uses full orthogonalization, as described in Algorithm 1, and the maximum eigenvalue of the Rayleigh matrix is computed with the command `eig`.

We consider several types of input matrices for which accurate maximum eigenvalue estimates are particularly difficult. The randomized block Krylov method is rotationally invariant and affine invariant, so there is no (theoretical) loss in working with a diagonal matrix. Of course, the Krylov method does not exploit the fact that the input matrix is diagonal.

Recall that an $n \times n$ matrix $\boldsymbol{W} \in \mathbb{R}^{n \times n}$ from the Gaussian Orthogonal Ensemble (GOE) is obtained by extracting the symmetric part of a Gaussian matrix:

$$\boldsymbol{W} = \tfrac{1}{2}(\boldsymbol{G} + \boldsymbol{G}^*) \in \mathbb{R}^{n \times n} \quad \text{where } \boldsymbol{G} \in \mathbb{R}^{n \times n} \text{ has independent standard normal entries.}$$

The large eigenvalues of a GOE matrix cluster together, and the spectral gap becomes increasingly small as the dimension increases. The spectrum has essentially no decay.

- To study the behavior of the Krylov method *without* a spectral gap, we draw a realization of a GOE matrix, diagonalize it using `eig`, and make an affine transformation so that its extreme eigenvalues are 0 and 1. Abusing terminology, we also refer to this diagonal input matrix as a GOE matrix.
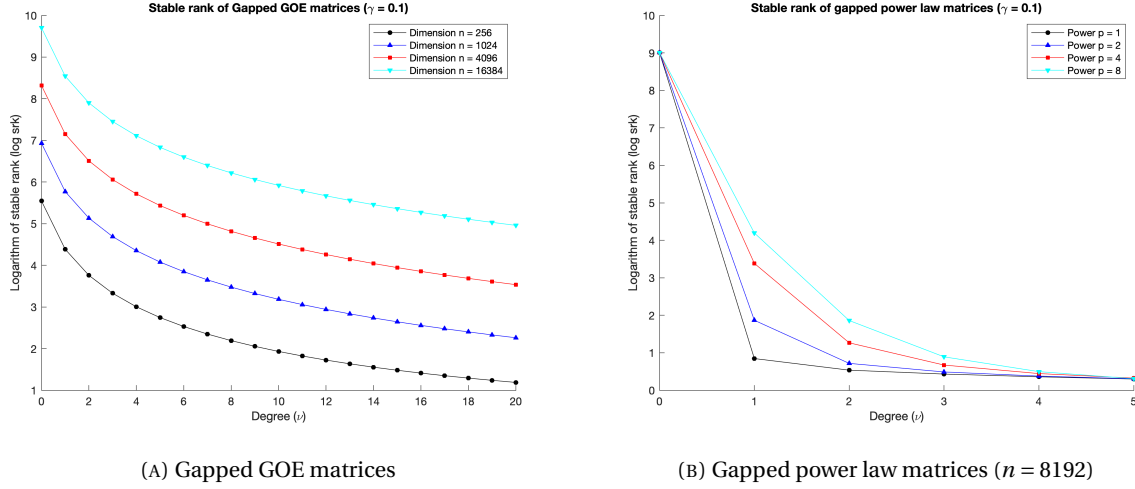
(A) Gapped GOE matrices

(B) Gapped power law matrices ($n = 8192$)

FIGURE 1: **Logarithm of stable rank.** These graphs display the natural logarithm of the stable rank function, $\mathrm{srk}(\nu)$, for two types of matrices. For gapped GOE matrices, the data series illustrate the effect of increasing the dimension ($n$). For gapped power law matrices, the series illustrate the effect of decreasing the rate of tail decay ($p$).

- To study the behavior of the Krylov method *with* a spectral gap $\gamma$, we take the diagonalized GOE matrix and increase the largest eigenvalue until the gap is $\gamma$. We refer to this model as a *gapped GOE matrix*.
- To understand how tail decay affects the convergence of the Krylov method, we consider *gapped power law* matrices. For a dimension $n$, power $p > 0$, and gap $\gamma \in [0, 1)$, the matrix is diagonal with entries
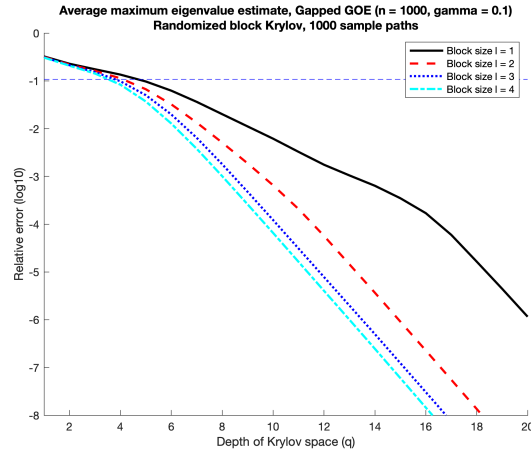
$$a_1 = 1 + \frac{\gamma}{1 - \gamma} \quad \text{and} \quad a_i = (i - 1)^{-1/p} \quad \text{for } i = 2, \ldots, n.$$

See Figure 1 for a depiction of the stable rank function for these types of input matrices. For the random matrix models, we draw a single realization of the input matrix and then fix it once and for all. The variability in the experiments derives from the random draw of the test matrix.

## 2.2. **Sample Paths: The Role of Block Size.**
First, we explore the empirical probability that the random-ized block Krylov method can achieve a sufficiently small error.

The first set of experiments focuses on the behavior of the block Krylov method for a matrix with a substantial spectral gap. For a fixed $1000 \times 1000$ gapped GOE matrix with $\gamma = 0.1$, Figure 2 illustrates sample paths of the relative error in estimating the maximum eigenvalue as a function of the total depth $q$ of the Krylov space. We compare the performance as the block size $\ell$ varies. Here is a summary of our observations:

- After a burn in period of about five steps, the error begins to decay exponentially, as described in Theorem 1.8. For $\ell = 4$, we can graphically estimate that the decay rate is about $e^{-1.38\,q}$, while the theory predicts a decay rate at least $e^{-1.26\,q}$.
- For block size $\ell = 1$, the average error decays at roughly *half* the rate achieved for block size $\ell \geq 2$. For $\ell = 2$, the average error has a distinctive profile; for $\ell = 3, 4$, the error is qualitatively similar. These phenomena match the bounds in Theorem 1.8(2).
- The sample paths give a clear picture of how the errors typically evolve. Independent of the block size, most of the sample paths decay at the same rate. As the block size increases, there is a slight reduction in the error (seen as a shift to the left), but this improvement is both small and diminishing.

(A) Average error across block size



(B) Block size $\ell = 1$



(C) Block size $\ell = 2$



(D) Block size $\ell = 3$



(E) Block size $\ell = 4$

FIGURE 2: **Sample paths, with spectral gap.** The input matrix is a fixed $1000 \times 1000$ GOE matrix with an artificial spectral gap $\gamma = 0.1$ (dashed blue line). The bottom four panels show the trajectory of the relative error (translucent hairlines) in estimating the maximum eigenvalue via randomized block Krylov for 1000 random test matrices with block size $\ell$ as a function of the depth $q$. The average error (**not** average log-error) is marked with a heavy orange line. The top panel compares the average error across block size. See Section 2.2.

- The impact of the block size becomes evident when we look at the spread of the sample paths. As the block size increases, the error varies much less, and this effect is sharpened by increasing the block size. The apparent reason is that the Krylov method can *misconverge*: it locks onto the *second largest* eigenvalue (dashed blue line), and it may not find a component of the maximum eigenvector for a significant number of iterations. For block size $\ell = 1$, this tendency is strong enough to destroy the rate of average convergence. For larger block sizes, the likelihood and duration of misconvergence both decrease.

The second set of experiments addresses the performance of the block Krylov method for a matrix with a small spectral gap. For a fixed $1000 \times 1000$ GOE matrix (with $\gamma \approx 0.0073$), Figure 3 shows sample paths of the relative error in estimating the maximum eigenvalue as a function of the total depth $q$ of the Krylov space and the block size $\ell$. Let us highlight a few observations.

- Initially, the average error decays roughly as $q^{-2}$, as predicted by Theorem 1.6. Once the error is sufficiently small (around $10^{-3}$), the decay becomes (super)exponential as indicated by Theorem 1.8.
- In the polynomial decay regime ($q \lesssim 20$), the block size affects the average error weakly, as suggested by Theorem 1.6. In the exponential decay regime, the block size plays a much more visible role, but the theory does not fully capture this effect.
- Regardless of the block size, the bulk of the sample paths decay more quickly than the average error. The block size has less of an effect on the typical error than on the average error. Nevertheless, for larger block size, we quickly achieve more digits of accuracy.
- As in the first experiment, the block size has a major impact on the variability of the error. As we increase the block size, the sample paths start to cluster sharply, so the typical error and the average error align with each other. Misconvergence is also visible here, and it can be mitigated by increasing the block size.

### 2.3. Burn-in: The Role of Tail Content.
Next, we examine how the tail content affects the burn-in period for the randomized block Krylov method.

Let us consider how the burn-in increases with the dimension of a matrix with limited spectral decay. For $n \times n$ gapped GOE matrices with $\gamma = 0.1$, Figure 4a shows how the average error evolves as a function of the depth $q$ of the Krylov space with block size $\ell = 2$. The main point is that the error initially stagnates before starting to decay exponentially. The stagnation period increases in proportion to the logarithm of the dimension $n$. As we saw in Figure 1a, this same behavior is visible in the size of $\log \mathrm{srk}(v)$ for small values of $v$.

Now, we look at how the burn-in depends on the rate of decay of the eigenvalues of a matrix. For an $8192 \times 8192$ gapped power law matrix with $\gamma = 0.1$, Figure 4b charts how the average error decays as a function of the depth $q$ for block size $\ell = 2$. We observe that the error initially decays exponentially at a slow rate. After a burn-in period, the error begins to decline at a faster exponential rate. The length of the initial trajectory increases in proportion with the logarithm of the order $p$ of the power law. From Figure 1b, we detect that this speciation reflects the spread of $\log \mathrm{srk}(v)$ for small values of $v$.

### 2.4. Conclusions.
The numerical experiments presented in this section confirm many of our theoretical predictions about the behavior of randomized block Krylov methods for estimating the maximum eigenvalue of a symmetric matrix. Increasing the block size results in a dramatic reduction in the probability of committing a large error in the estimate. The underlying mechanism is that simple Krylov method ($\ell = 1$) is much more likely to misconverge than a block method (even with $\ell = 2, 3$). The effect is so significant that block Krylov methods can converge, on average, at twice the rate of a simple Krylov method. These benefits counterbalance the increased computational cost of a block method.

These facts have implications for design and implementation of practical algorithms. If the arithmetic cost is the driving concern, then the simple randomized Krylov method typically uses matrix–vector multiplications more efficiently than the block methods. But we have also seen that block methods can
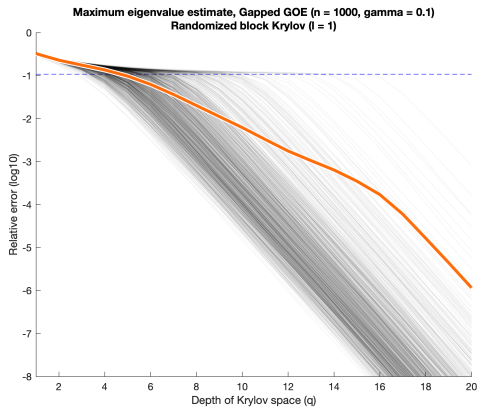
(A) Average error across block size



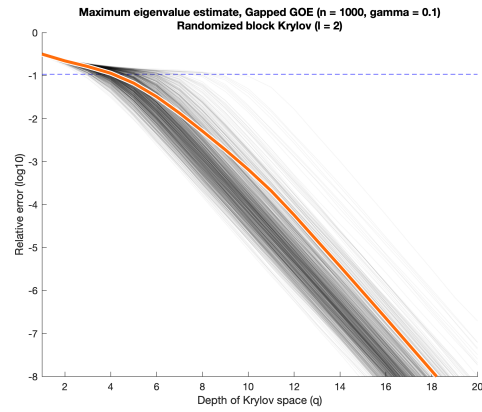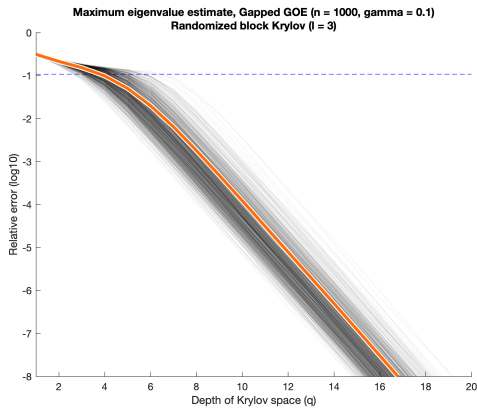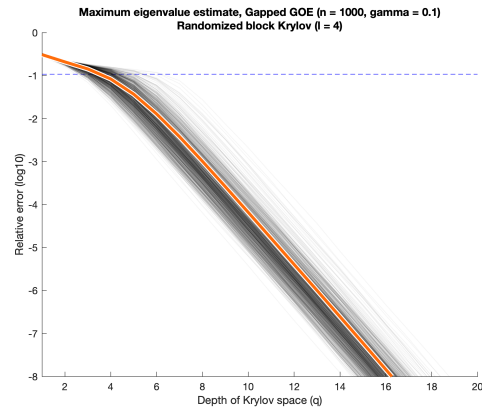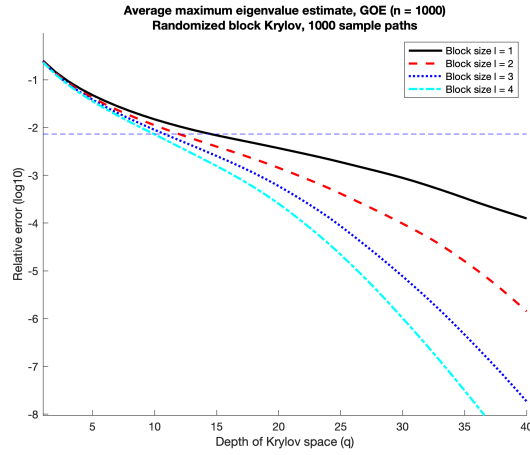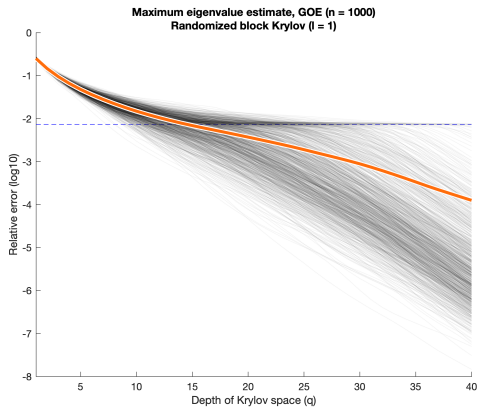(B) Block size $\ell = 1$



(C) Block size $\ell = 2$



(D) Block size $\ell = 3$



(E) Block size $\ell = 4$

FIGURE 3: **Sample paths, no spectral gap.** The input matrix is a fixed $1000 \times 1000$ GOE matrix with spectral gap $\gamma \approx 0.0073$ (dashed blue line). The bottom four panels illustrate the trajectory of the relative error (translucent hairlines) in estimating the maximum eigenvalue via the randomized block Krylov method for 1000 random test matrices with block size $\ell$ as a function of the depth $q$. The average error (**not** average log-error) is marked with a heavy orange line. The top panel compares the average error across block size. See Section 2.2.

(A) Role of dimension                    (B) Role of tail decay

FIGURE 4: **Burn-in as a function of tail content.** These plots show the average error in estimating the maximum eigenvalue as a function of depth $q$ with block size $\ell = 2$. In the left panel, we consider gapped GOE matrices of increasing dimension $n$. In the right panel, we consider gapped power law matrices with decreasing rate $p$ of tail decay.

achieve a specific target for the error (averaged over test matrices) using a similar number of matrix–vector multiplies as the simple method. In fact, in modern computing environments, the actual cost (time, communication, energy) of multiplying a matrix by several vectors may be equivalent to a single matrix–vector multiplication, in which case the block methods are the obvious choice. Finally, when it is critical to limit the probability of failure, then the block methods are clearly superior.

## 3. HISTORY, RELATED WORK, AND CONTRIBUTIONS

Krylov subspace methods are a wide class of algorithms that use matrix–vector products ("Krylov information") to compute eigenvalues and eigenvectors and to solve linear systems. These methods are especially appropriate in situations where we can only interact with a matrix through its action on vectors. In this treatment, we only discuss Krylov methods for spectral computations. Some of the basic algorithms in this area are 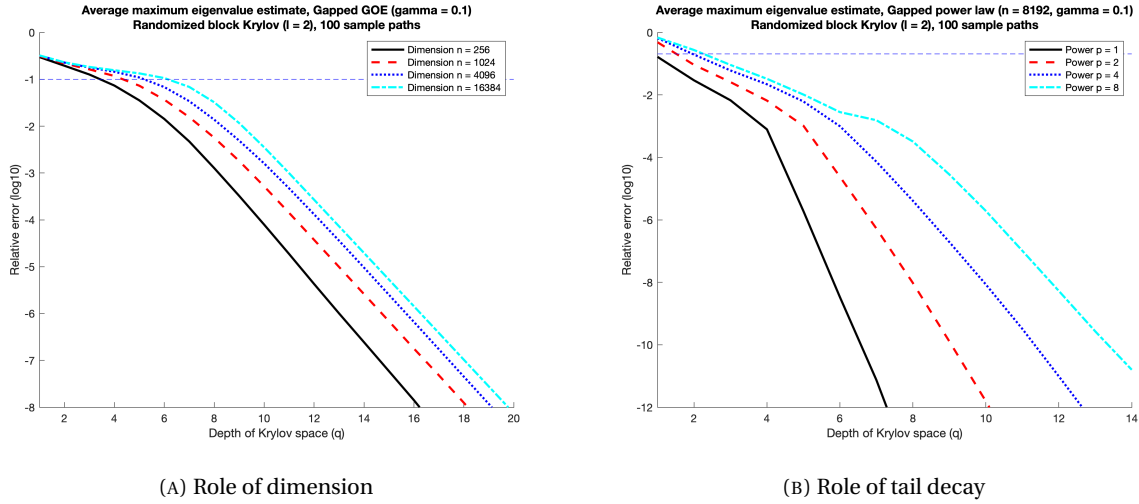the power method, the inverse power method, subspace iteration, the Lanczos method, the block Lanczos method, and the band Lanczos method. See the books [Par98, BDD+00, Saa11, GVL13] for more background and details.

3.1. **Simple Krylov Methods.** Simple Krylov methods are algorithms based on a Krylov subspace $K_q(A; b)$ constructed from a single starting vector $b$. That is, the block size $\ell = 1$.

The power method, which dates to the 19th century, is probably the earliest algorithm that relies on Krylov information to compute eigenvalues and eigenvectors of a symmetric matrix. The power method is degenerate in the sense that it only keeps the highest-order term in the Krylov subspace.

In the late 1940s, Lanczos [Lan50] developed a sophisticated Krylov subspace method for solving the symmetric eigenvalue problem. (More precisely, the Lanczos method uses a three-term recurrence to compute a basis for the Krylov subspace so that the compression of the input matrix to the Krylov subspace is tridiagonal.) In exact arithmetic, the Lanczos estimate of the maximum eigenvalue of a symmetric matrix coincides with $\xi_{\max}(A; b; q)$ for a fixed vector $b$. On the other hand, the Lanczos method has complicated behavior in finite-precision arithmetic; for example, see Meurant's book [Meu06].

The first analysis of the Lanczos method with a deterministic starting vector $b$ dates to the work of Lanczos [Lan50]. Kaniel, Paige, and Saad also made major theoretical contributions in the 1970s and 1980s; see [Par98, Saa11] for details and references. In the 1980s, Nemirovsky, Yudin, and Chou showed that Krylov

subspace methods are the optimal deterministic algorithms for solving the symmetric eigenvalue problem, assuming we only have access to the matrix via matrix–vector multiplication; see [NY83, Cho87, Nem91].

The burn-in period for Krylov methods has been observed in many previous works, including [vdSvdV86, LS05, BES05]. The length of the burn-in period depends on the proportion of the energy in the test vector that is captured by the true invariant subspace.

The main contribution of the paper [vdSvdV86] is the observation that Krylov methods may exhibit superlinear convergence. The explanation for this phenomenon is that the optimal polynomial can create an artificial spectral gap by annihilating undesirable eigenvalues. This kind of behavior is visible in Figure 3.

### 3.2. **Random Starting Vectors.**

Practitioners have often suggested using randomized variants of Krylov subspace methods. That is, the starting vector $b$ is chosen at random. Historically, randomness was just used to avoid the situation where the starting vector $b$ is orthogonal to the leading invariant subspace of the matrix.

Later, deeper justifications for random starting vectors appeared. The first substantive theoretical analysis of a randomized Krylov method appears in Dixon's paper [Dix83] on the power method with a random starting vector. We believe that this is the first paper to recognize that Krylov methods can be successful without the presence of a spectral gap.

In 1992, Kuczyński & Woźniakowski published an analysis [KW92] of the Lanczos method with a random starting vector. Their work highlighted the benefits of randomization, and it provided a clear explanation of the advantages of using full Krylov information instead of the power method. See the papers [KW94, DCM97, LW98] for further work in this direction.

The recent papers [SAR17, SEAR18] contain lower bounds on the performance of randomized algorithms for the symmetric eigenvalue problem that use Krylov information.

### 3.3. **Block Krylov Methods.**

Block Krylov subspace methods use multiple starting vectors to generate the Krylov subspace, instead of just one. In other words, the algorithms form a Krylov subspace $K_q(A; B)$, where $B$ is a matrix. These methods were developed in the late 1960s and 1970s in an effort to resolve multiple eigenvalues more reliably. The block analog of the power method is called subspace iteration; see the books [Par98, Saa11] for discussion.

There are also block versions of the Lanczos method, which were developed by Cullum & Donath [CD74] and Golub & Underwood [GU77]. (More precisely, the block Lanczos method uses a recurrence to compute a basis for the block Krylov subspace so that the compression of the input matrix to the block Krylov subspace is block tridiagonal.) In exact arithmetic, the block Lanczos estimate of the maximum eigenvalue of a symmetric matrix coincides with $\xi_{\max}(A; B; q)$ for a fixed matrix $B$.

Most of the early work on block Krylov subspace methods focuses on the case where the block size $\ell$ is small, while the depth $q$ of the Krylov space is moderately large. This can lead to significant problems with numerical stability, especially in the case where we use a recurrence to perform orthogonalization. As a consequence, full orthogonalization is usually recommended. Furthermore, most of the existing analysis of block Krylov methods is deterministic; for example, see [Saa80, LZ15].

### 3.4. **Randomized Block Krylov Methods.**

Over the last decade, randomized block Krylov subspace methods have emerged as a powerful tool for spectral computations on large matrices. These algorithms use a Krylov subspace $K_q(A; B)$ generated by a *random* test matrix $B$.

In contrast with earlier block Krylov algorithms, contemporary methods use a much larger block size $\ell$ and a much smaller depth $q$. Furthermore, the randomness plays a central role in supporting performance guarantees for the algorithms.

Most of the recent literature concerns the problem of computing a low-rank approximation to a large matrix, rather than estimating eigenvalues or invariant subspaces. Some of the initial work on randomized algorithms for matrix approximation appears in [PRTV98, FKV98, DKM06, MRT11]. Randomized subspace iteration was proposed in [RST09] and further developed in [HMT11]. Randomized block Krylov methods

that use the full block Krylov subspace were proposed in the papers [RST09, HMST11]; see also [DIKMI18]. See [HMT11] for more background and history.

There is some theoretical and empirical evidence [HMST11, MM15] that randomized block Krylov methods can produce low-rank matrix approximations with higher accuracy and with less computation than randomized subspace iteration.

### 3.5. **Analysis of Randomized Block Krylov Methods.**

There is a growing body of literature that develops theoretical performance guarantees for randomized block Krylov methods. The papers [RST09, HMT11, Gu15, MM15] contain theoretical analyses of randomized subspace iteration. The papers [MM15, WZZ15, DIKMI18] contain theoretical analysis of randomized methods that use the full block Krylov space. These works all focus on low-rank matrix approximation.

### 3.6. **Contemporary Work.**

After the initial draft of this paper and its companion [Tro18] was completed, several other research groups released new work on randomized (block) Krylov methods. The paper [DI19] shows that randomized block Lanczos can produce excellent low-rank matrix approximations, even without the presence of a spectral gap. The paper [YGL18] demonstrates that randomized block Lanczos gives excellent estimates for the singular values of a general matrix with spectral decay.

### 3.7. **Contributions.**

We set out to develop highly refined bounds for the behavior of randomized block Krylov methods that use the full Krylov subspace. Our aim is to present useful and informative results in the spirit of Saad [Saa80], Kucziński & Woźniakowski [KW92], and Halko et al. [HMT11]. Our research has a number of specific benefits over prior work.

- We have shown that randomized block Krylov methods have exceptional performance for matrices with spectral decay. In fact, for matrices with polynomial spectral decay, we can often obtain accurate estimates even when the block Krylov subspace has *constant* depth.
- We have obtained detailed information about the role of the block size $\ell$. In particular, increasing the block size drives down failure probabilities exponentially fast.
- The companion paper [Tro18] gives the first results on the performance of randomized block Krylov methods for the symmetric eigenvalue problem.
- Our bounds have explicit and modest constants, which gives the bounds some predictive power.

We hope that these results help clarify the benefits of randomized block Krylov methods. We also hope that our work encourages researchers to develop new implementations of these algorithms that fully exploit contemporary computer architectures.

## 4. PRELIMINARIES

Before we begin the proofs of the main results, Theorems 1.6 and 1.8, we present some background information. In Section 4.1, we justify the claim that the test matrix should have a uniformly random range. In Section 4.2, we introduce the Chebyshev polynomials of the first and second type, and we develop the properties that we need to support our arguments. Expert readers may wish to skip to Section 5 for proofs of the main results.

### 4.1. **Rotationally Invariant Distributions.**

We wish to estimate spectral properties of a matrix using randomized block Krylov information. In particular, we aim to establish probabilistic upper bounds on the error in these spectral estimates for any input matrix. This section contains a general argument that clarifies why we ought to use a random test matrix with a rotationally invariant range in these applications.

**Proposition 4.1** (Uniformly Random Range). *Consider any bivariate function $f : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times \ell} \to \mathbb{R}$ that is orthogonally invariant:*

$$f(\boldsymbol{A}; \boldsymbol{B}) = f(\boldsymbol{U}\boldsymbol{A}\boldsymbol{U}^*; \boldsymbol{U}\boldsymbol{B}) \quad \textit{for each orthogonal } \boldsymbol{U} \in \mathbb{R}^{n \times n}.$$

*Fix a symmetric $n \times n$ matrix $\boldsymbol{\Lambda}$, and consider the orthogonal orbit*

$$\mathscr{A} := \mathscr{A}(\boldsymbol{\Lambda}) := \left\{ \boldsymbol{U}\boldsymbol{\Lambda}\boldsymbol{U}^* : \boldsymbol{U} \in \mathbb{R}^{n \times n} \text{ orthogonal} \right\}.$$

*Let $B \in \mathbb{R}^{n \times \ell}$ be a random matrix. Let $V \in \mathbb{R}^{n \times n}$ be a uniformly random orthogonal matrix, drawn independently from $B$. Then*

$$\max_{A \in \mathscr{A}} \mathbb{E}_{V,B} f(A; VB) \leq \max_{A \in \mathscr{A}} \mathbb{E}_B f(A; B).$$

*Proof.* By orthogonal invariance of $f$,

$$\max_{A \in \mathscr{A}} \mathbb{E}_{V,B} f(A; VB) = \max_{A \in \mathscr{A}} \mathbb{E}_{V,B} f(V^* AV; B)$$

$$\leq \mathbb{E}_V \max_{A \in \mathscr{A}} \mathbb{E}_B f(V^* AV; B) = \max_{A \in \mathscr{A}} \mathbb{E}_B f(A; B).$$

The inequality is Jensen's, and the last identity follows from the definition of the class $\mathscr{A}$. $\qquad\square$

Now, consider the problem of estimating the maximum eigenvalue of the worst matrix with eigenvalue spectrum $\Lambda$ using a block Krylov method. Define the orthogonally invariant[1] function $f(A; B) = \mathrm{err}(\xi_{\max}(A; B; q))$. For any distribution on $B$ and for a uniformly random orthogonal matrix $V$, Proposition 4.1 states that

$$\max_{A \in \mathscr{A}} \mathbb{E}\, \mathrm{err}(\xi_{\max}(A; VB; q)) \leq \max_{A \in \mathscr{A}} \mathbb{E}\, \mathrm{err}(\xi_{\max}(A; B; q)).$$

That is, the random test matrix $VB$ is better than the test matrix $B$ if we want to minimize the worst-case expectation of the error; the same kind of bound holds for tail probabilities. We surmise that the test matrix $B$ should have a uniformly random range. Moreover, because of the co-range invariance (1.6), we can select *any* distribution with uniformly random range, such as the standard normal matrix $\Omega$.

### 4.2. **Chebyshev Polynomials.**

In view of (1.3), we can interpret Krylov subspace methods as algorithms that compute best polynomial approximations. The analysis of Krylov methods often involves a careful choice of a specific polynomial that witnesses the behavior of the algorithm. Chebyshev polynomials arise naturally in this connection because they are the solutions to minimax approximation problems.

This section contains the definitions of the Chebyshev polynomials of the first and second type, and it derives some key properties of these polynomials. We also construct the specific polynomials that arise in our analysis. Some of this material is drawn from the paper [KW92] by Kuczyński & Woźniakowski. For general information about Chebyshev polynomials, we refer the reader to [AS64, OLBC10].

4.2.1. *Chebyshev Polynomials of the First Kind.* We can define the Chebyshev polynomials of the first kind via the formula

$$T_p(s) := \frac{1}{2}\left[\left(s + \sqrt{s^2 - 1}\right)^p + \left(s - \sqrt{s^2 - 1}\right)^p\right] \quad \text{for } s \in \mathbb{R} \text{ and } p \in \mathbb{Z}_+. \tag{4.1}$$

Using the binomial theorem, it is easy to check that this expression coincides with a polynomial of degree $p$ with real coefficients.

We require two properties of the Chebyshev polynomial $T_p$. First, it satisfies a uniform bound on the unit interval:

$$\left|T_p(s)\right| \leq 1 \quad \text{for } |s| \leq 1. \tag{4.2}$$

Indeed, it is well-known that $2^{-p} T_p$ is the unique monic polynomial of degree $p$ with the least maximum value on the interval $[-1, +1]$. The simpler result (4.2) is an immediate consequence of the representation

$$T_p(s) = \cos\left(p \cos^{-1}(s)\right) \quad \text{for } |s| \leq 1.$$

The latter formula follows from (4.1) after we apply de Moivre's theorem for complex exponentiation.

Second, the Chebyshev polynomial grows quickly outside of the unit interval:

$$T_p\left(\frac{1+s}{1-s}\right) \geq \frac{1}{2}\left(\frac{1+\sqrt{s}}{1-\sqrt{s}}\right)^p \quad \text{for } 0 \leq s < 1. \tag{4.3}$$

This estimate is a direct consequence of the definition (4.1).

---

[1]The relative error (1.12) is orthogonally invariant because the eigenvalues of $A$ are orthogonally invariant and (1.7) states that the eigenvalue estimate (1.4) is also orthogonally invariant.

4.2.2. *The Attenuation Factor.* Let $\beta \in [0,1]$ be a parameter, and define the quantity

$$\delta := \delta(\beta) = \frac{1 - \sqrt{1-\beta}}{1 + \sqrt{1-\beta}}. \tag{4.4}$$

This definition is closely connected with the growth properties of $T_p$. We can bound the attenuation factor in two ways:

$$\delta \le e^{-2\sqrt{1-\beta}} \quad \text{and} \quad \delta \le \beta \cdot 2^{-2\sqrt{1-\beta}}. \tag{4.5}$$

These numerical inequalities can be justified using basic calculus. The first is very accurate for $\beta \approx 1$, while the second is better across the full range $\beta \in [0,1]$.

4.2.3. *First Polynomial Construction.* Choose a nonnegative integer partition $q = q_1 + q_2$. Consider the polynomial

$$\varphi_{\beta,q_1,q_2}(s) := \frac{s^{q_1} T_{q_2}((2/\beta)s - 1)}{T_{q_2}((2/\beta) - 1)} \quad \text{for } s \in \mathbb{R}.$$

The polynomial $\varphi_{\beta,q_1,q_2}$ has degree $q$, and it is normalized so that it takes the value one at $s = 1$. It holds that

$$\varphi_{\beta,q_1,q_2}^2(s) \le \frac{s^{2q_1}}{T_{q_2}^2((2/\beta) - 1)} \le 4s^{2q_1} \left( \frac{1 - \sqrt{1-\beta}}{1 + \sqrt{1-\beta}} \right)^{2q_2} = 4s^{2q_1}\delta^{2q_2} \quad \text{for } 0 \le s \le \beta. \tag{4.6}$$

The first inequality follows from (4.2), and the second follows from (4.3). Last, we instate the definition (4.4).

**Remark 4.2** (The Monomial)**.** In contrast to [KW92] and other prior work, we use products of Chebyshev polynomials with low-degree monomials. This seemingly minor change leads to results phrased in terms of the stable rank, rather than the ambient dimension. This is equivalent to a bound on $q_1$ iterations of the subspace iteration method, followed by $q_2$ iterations of the block Krylov method.

4.2.4. *Chebyshev Polynomials of the Second Kind.* We can define the Chebyshev polynomials of the second kind via the formula

$$U_p(s) := \frac{1}{2\sqrt{s^2 - 1}} \left[ \left( s + \sqrt{s^2 - 1} \right)^{p+1} - \left( s - \sqrt{s^2 - 1} \right)^{p+1} \right] \quad \text{for } s \in \mathbb{R} \text{ and } p \in \mathbb{Z}_+. \tag{4.7}$$

Using the binomial theorem, it is easy to check that this expression coincides with a polynomial of degree $p$ with real coefficients. Moreover, when $p$ is an even number, the polynomial $U_p$ is an even function.

We require two properties of the Chebyshev polynomial $U_p$. First, it satisfies a weighted uniform bound on the unit interval:

$$\left| \sqrt{1 - s^2}\, U_p(s) \right| \le 1 \quad \text{for } |s| \le 1. \tag{4.8}$$

In fact, $2^{-p}U_p$ is the unique monic polynomial that minimizes the maximum value of the left-hand side of (4.8) over the interval $[-1, +1]$; see [KW92, Eqn. (23) et seq.]. The simpler result (4.8) is an immediate consequence of the representation

$$U_p(s) = \frac{\sin((p+1)\cos^{-1}(s))}{\sqrt{1 - s^2}} \quad \text{for } |s| \le 1.$$

The latter formula follows from (4.7) after we apply de Moivre's theorem for complex exponentiation.

Second, we can evaluate the polynomial $U_{2p}$ at a specific point:

$$U_{2p}^2\left(\sqrt{1/\beta}\right) = \frac{\beta\left(1 - \delta^{2p+1}\right)^2}{4(1 - \beta)\delta^{2p+1}} \quad \text{where } 0 < \beta \le 1. \tag{4.9}$$

We defined $\delta = \delta(\beta)$ above in (4.4). The formula (4.9) is a direct—but unpleasant—consequence of the definition (4.7).

4.2.5. *Second Polynomial Construction.* As before, introduce a parameter $\beta \in [0,1]$. Choose a nonnegative integer partition $q = q_1 + q_2$. Consider the polynomial

$$\psi_{\beta,q_1,q_2}(s) := \frac{s^{q_1} U_{2q_2}\left(\sqrt{s/\beta}\right)}{U_{2q_2}\left(\sqrt{1/\beta}\right)} \quad \text{for } s \in \mathbb{R}.$$

Since $U_{2q_2}$ is an even polynomial, this expression defines a polynomial $\psi_{\beta,q_1,q_2}$ with degree $q$ and normalized to take the value one at $s = 1$. We have the bound

$$(\beta - s)\,\psi^2_{\beta,q_1,q_2}(s) \le \frac{s^{2q_1}\beta}{U^2_{2q_2}(\sqrt{1/\beta})} = \frac{4(1-\beta)s^{2q_1}\delta^{2q_2+1}}{\left(1 - \delta^{2q_2+1}\right)^2} \quad \text{for } 0 \le s \le \beta. \tag{4.10}$$

The inequality (4.10) follows from (4.8), and the equality follows from (4.9). Last, we instate the definition (4.4). For $s > \beta$, the polynomial grows very rapidly.

## 5. The Error in the Block Krylov Subspace Method

In this section, we initiate the proof of Theorems 1.6 and 1.8. Along the way, we establish Proposition 1.5. First, we show how to replace the block Krylov subspace by a simple Krylov subspace (with block size one). Afterward, we develop an explicit representation for the error in the eigenvalue estimate derived from the simple Krylov subspace. Finally, we explain how to construct the simple Krylov subspace so that we preserve the benefits of computing a block Krylov subspace.

The ideas in this section are drawn from several sources. The strategy of reducing a block Krylov subspace to a simple Krylov subspace already appears in [Saa80], but we use a different technique that is adapted from [HMT11]. The kind of analysis we perform for the simple Krylov method is standard; we have closely followed the presentation in [KW92].

5.1. **Simplifications.** Suppose that the input matrix $\boldsymbol{A}$ is a multiple of the identity matrix. From the definitions (1.1), (1.2), and (1.4), it is straightforward to check that the eigenvalue estimate $\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q) = \lambda_{\max}(\boldsymbol{A})$ with probability one for each $q \ge 0$. Therefore, we may as well assume that $\boldsymbol{A}$ is not a multiple of the identity.

In view of (1.9), we may also assume that the input matrix is diagonal with weakly decreasing entries:

$$\boldsymbol{A} = \operatorname{diag}(a_1, a_2, \ldots, a_n) \quad \text{where} \quad a_1 \ge a_2 \ge \cdots \ge a_n. \tag{5.1}$$

Since $\boldsymbol{A}$ has at least two distinct eigenvalues, we may normalize the extreme eigenvalues of $\boldsymbol{A}$:

$$a_1 = 1 \quad \text{and} \quad a_n = 0. \tag{5.2}$$

The main results are all stated in terms of affine invariant quantities, so we have not lost any generality. These choices help to streamline the proof.

5.2. **Block Krylov Subspaces and Simple Krylov Subspaces.** The first key step in the argument is to reduce the block Krylov subspace to a Krylov subspace with block size one. This idea allows us to avoid any computations involving matrices. To that end, recall that the block Krylov subspace takes the form

$$K_q(\boldsymbol{A}; \boldsymbol{\Omega}) = \operatorname{range}\begin{bmatrix} \boldsymbol{\Omega} & \boldsymbol{A}\boldsymbol{\Omega} & \boldsymbol{A}^2\boldsymbol{\Omega} & \ldots & \boldsymbol{A}^q\boldsymbol{\Omega} \end{bmatrix}.$$

In particular, for any vector $\boldsymbol{x} \in \operatorname{range}(\boldsymbol{\Omega})$,

$$K_q(\boldsymbol{A}; \boldsymbol{x}) = \operatorname{range}\begin{bmatrix} \boldsymbol{x} & \boldsymbol{A}\boldsymbol{x} & \boldsymbol{A}^2\boldsymbol{x} & \ldots & \boldsymbol{A}^q\boldsymbol{x} \end{bmatrix} \subset K_q(\boldsymbol{A}; \boldsymbol{\Omega}).$$

Later, we will make a careful choice of the vector $\boldsymbol{x}$ so that we do not abandon the benefits of computing the block Krylov subspace.

5.3. **Representation of the Error Using Polynomials.** The next step in the argument is to exploit the close relationship between Krylov subspaces and polynomial filtering to obtain an explicit representation of the error in the eigenvalue estimate. Using (1.3), we may rewrite the last display in the form

$$K_q(\boldsymbol{A}; \boldsymbol{x}) = \operatorname{span}\{\varphi(\boldsymbol{A})\boldsymbol{x} : \varphi \in \mathscr{P}_q\} \subset K_q(\boldsymbol{A}; \boldsymbol{\Omega}).$$

As a consequence of this containment,

$$\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q) = \max_{\boldsymbol{v} \in K_q(\boldsymbol{A}; \boldsymbol{\Omega})} \frac{\boldsymbol{v}^* \boldsymbol{A} \boldsymbol{v}}{\boldsymbol{v}^* \boldsymbol{v}} \geq \max_{\boldsymbol{v} \in K_q(\boldsymbol{A}; \boldsymbol{x})} \frac{\boldsymbol{v}^* \boldsymbol{A} \boldsymbol{v}}{\boldsymbol{v}^* \boldsymbol{v}} = \max_{\varphi \in \mathscr{P}_q} \frac{(\varphi(\boldsymbol{A})\boldsymbol{x})^* \boldsymbol{A} (\varphi(\boldsymbol{A})\boldsymbol{x})}{(\varphi(\boldsymbol{A})\boldsymbol{x})^* (\varphi(\boldsymbol{A})\boldsymbol{x})}.$$

Owing to the normalization (5.2), the relative error (1.12) in the eigenvalue estimate satisfies

$$\operatorname{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) = 1 - \xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q) \leq \min_{\varphi \in \mathscr{P}_q} \frac{(\varphi(\boldsymbol{A})\boldsymbol{x})^* (\boldsymbol{I} - \boldsymbol{A}) (\varphi(\boldsymbol{A})\boldsymbol{x})}{(\varphi(\boldsymbol{A})\boldsymbol{x})^* (\varphi(\boldsymbol{A})\boldsymbol{x})}.$$

The fraction is continuous on the set of polynomials $\{\varphi \in \mathscr{P}_q : \varphi(\boldsymbol{A})\boldsymbol{x} \neq \boldsymbol{0}\}$, and it is invariant under scaling of the polynomial $\varphi$. Thus, assuming that the first coordinate of $\boldsymbol{x}$ is nonzero, we may assume that $\varphi(1) \neq 0$ and then rescale so that $\varphi(1) = 1$. With the definition $\mathscr{P}_q(1) := \{\varphi \in \mathscr{P}_q : \varphi(1) = 1\}$, we arrive at

$$\operatorname{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \inf_{\varphi \in \mathscr{P}_q(1)} \frac{(\varphi(\boldsymbol{A})\boldsymbol{x})^* (\boldsymbol{I} - \boldsymbol{A}) (\varphi(\boldsymbol{A})\boldsymbol{x})}{(\varphi(\boldsymbol{A})\boldsymbol{x})^* (\varphi(\boldsymbol{A})\boldsymbol{x})}.$$

We remark that this inequality becomes an equality in the case where the block size $\ell = 1$.

Invoke (5.1) to rewrite this bound in terms of the eigenvalues of $\boldsymbol{A}$:

$$\operatorname{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \inf_{\varphi \in \mathscr{P}_q(1)} \frac{\sum_{i=1}^n X_i^2 \varphi^2(a_i)(1 - a_i)}{\sum_{i=1}^n X_i^2 \varphi^2(a_i)} = \inf_{\varphi \in \mathscr{P}_q(1)} \frac{\sum_{i>1} X_i^2 \varphi^2(a_i)(1 - a_i)}{X_1^2 + \sum_{i>1} X_i^2 \varphi^2(a_i)}. \tag{5.3}$$

We have introduced the coordinates of the vector $\boldsymbol{x} = (X_1, \dots, X_n)$, and we have used normalization (5.2) to simplify the expression.

**Remark 5.1** (Multiplicity)**.** When the maximum eigenvalue has multiplicity $m$ greater than one, we can group the copies of the eigenvalue together in the denominator to obtain the larger term $X_1^2 + \cdots + X_m^2$. Pursuing this argument, we see that the algorithm converges faster when $m > 1$.

5.4. **When the Matrix has Few Distinct Eigenvalues.** Suppose that the distinct eigenvalues of the input matrix $\boldsymbol{A}$ are $\mu_1 = 1$ and $\mu_2, \dots, \mu_r$ where $r \leq q + 1$. Consider the polynomial

$$\varphi_0(s) = \prod_{i=2}^r \frac{s - \mu_i}{1 - \mu_i} \in \mathscr{P}_{r-1}(1) \subset \mathscr{P}_q(1).$$

This polynomial annihilates each point in the spectrum, except for $\mu_1$. Setting $\varphi = \varphi_0$ in (5.3), we discover that $\operatorname{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) = 0$ with probability one. This fact appears as Proposition 1.5.

5.5. **Choosing the Simple Krylov Space.** The next step in the argument is to select a particular vector $\boldsymbol{x} \in \operatorname{range}(\boldsymbol{\Omega})$. Let $\boldsymbol{\omega}_1^* \in \mathbb{R}^\ell$ denote the first row of the matrix $\boldsymbol{\Omega}$. Set

$$\boldsymbol{x} = \frac{\boldsymbol{\Omega}\boldsymbol{\omega}_1}{\|\boldsymbol{\omega}_1\|} \in \operatorname{range}(\boldsymbol{\Omega}).$$

The rows of $\boldsymbol{\Omega}$ are statistically independent standard normal vectors, which are rotationally invariant. It follows that the entries of $\boldsymbol{x} = (X_1, \dots, X_n)$ are also statistically independent. Moreover,

$$X_1 \sim \text{CHI}(\ell) \quad \text{and} \quad X_i \sim \text{NORMAL}(0, 1) \quad \text{for } i > 1. \tag{5.4}$$

We write $\text{CHI}(\ell)$ for the chi distribution with $\ell$ degrees of freedom. This choice of $\boldsymbol{x}$ ensures that $X_1^2$ is large relative to the other $X_i^2$. In particular, $X_1 \neq 0$ with probability one.

In a concrete sense, we see that the block Krylov method is an artifice for increasing the apparent multiplicity of the maximum eigenvalue. The argument here is inspired by the analysis in Halko et al. [HMT11, Secs. 9, 10].

**Remark 5.2** (Complex Setting)**.** When the input matrix is Hermitian and the test matrix follows the complex standard normal distribution, we obtain the same formulas with $X_1 \sim \text{CHI}(2\ell)$ and $X_i \sim \text{NORMAL}_{\mathbb{C}}(0,1)$. In effect, the multiplicity of every eigenvalue is doubled.

## 6. Probabilistic Bounds for the Error

This section contains the proof of the probability bounds that appear in Theorem 1.6 and 1.8. The argument is based on the bound (5.3) for the error and the distributional properties (5.4) of the random vector $\boldsymbol{x}$.

The proof is inspired by the argument in Kuczyński & Woźniakowski [KW92], but our approach is technically easier. Indeed, they work with a random vector that is uniformly distributed on the Euclidean unit sphere, which leads to a difficult multivariate integration. In contrast, our random vector $\boldsymbol{x}$ has independent entries, which means that we only have to compute one-dimensional integrals.

### 6.1. A Bound for the Probability.

Let $\varepsilon \in (0,1)$ be an error tolerance. Our goal is to control the probability $P_\varepsilon$ that the relative error in the eigenvalue estimate (1.4) is at least $\varepsilon$. In other words, we wish to bound

$$P_\varepsilon := \mathbb{P}\left\{ \text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \geq \varepsilon \right\}.$$

In view of the upper bound (5.3) for the relative error, we obtain the estimate

$$P_\varepsilon \leq \mathbb{P}\left\{ \inf_{\varphi \in \mathscr{P}_q(1)} \frac{\sum_{i>1} X_i^2 \varphi^2(a_i)(1 - a_i)}{X_1^2 + \sum_{i>1} X_i^2 \varphi^2(a_i)} \geq \varepsilon \right\}.$$

Fix a polynomial $\varphi \in \mathscr{P}_q(1)$, to be determined later. Then rearrange the inequality in the event:

$$
\begin{aligned}
P_\varepsilon &\leq \mathbb{P}\left\{ \sum\nolimits_{i>1} X_i^2 \varphi^2(a_i)(1 - a_i) \geq \varepsilon X_1^2 + \varepsilon \sum\nolimits_{i>1} X_i^2 \varphi^2(a_i) \right\} \\
&= \mathbb{P}\left\{ -\varepsilon X_1^2 + \sum\nolimits_{i>1} X_i^2 \varphi^2(a_i)(1 - \varepsilon - a_i) \geq 0 \right\} \\
&\leq \mathbb{P}\left\{ -\varepsilon X_1^2 + \sum\nolimits_{a_i < 1-\varepsilon} X_i^2 \varphi^2(a_i)(1 - \varepsilon - a_i) \geq 0 \right\} \\
&=: \mathbb{P}\left\{ -\varepsilon X_1^2 + \sum\nolimits_{i \in I} c_i X_i^2 \geq 0 \right\}.
\end{aligned}
\tag{6.1}
$$

To reach the third line, we have dropped the nonpositive terms in the sum. Then we introduced the compact notation

$$I := \{ i : a_i < 1 - \varepsilon \} \quad \text{and} \quad c_i := \varphi^2(a_i)(1 - \varepsilon - a_i) > 0.$$

To continue the argument, we apply some elementary notions from the theory of concentration of measure.

### 6.2. The Laplace Transform Argument.

We invoke the Laplace transform method to convert the probability bound into an expectation bound. Introduce a parameter $\theta > 0$, to be chosen later. Continuing from (6.1), we write the probability as the expectation of a 0–1 indicator function:

$$
\begin{aligned}
P_\varepsilon &\leq \mathbb{E}\mathbb{1}\left\{ -\varepsilon X_1^2 + \sum\nolimits_{i \in I} c_i X_i^2 \geq 0 \right\} \\
&\leq \mathbb{E}\exp\left( -\theta \varepsilon X_1^2 + \sum\nolimits_{i \in I} \theta c_i X_i^2 \right).
\end{aligned}
$$

To reach the second line, we bound the indicator $\mathbb{1}\{s \geq 0\}$ above by the function $s \mapsto e^{\theta s}$. Write the exponential as a product, and invoke the independence of the family $\{X_i\}$ to obtain

$$P_\varepsilon \leq \mathbb{E}\left[ e^{-\theta \varepsilon X_1^2} \prod\nolimits_{i \in I} e^{\theta c_i X_i^2} \right] = \left( \mathbb{E} e^{-\theta \varepsilon X_1^2} \right)\left( \prod\nolimits_{i \in I} \mathbb{E} e^{\theta c_i X_i^2} \right).$$

The distributional property (5.4) implies that $X_1^2$ is a chi-squared variable with $\ell$ degrees of freedom, while $X_i^2$ is a chi-squared variable with one degree of freedom. Computing the remaining expectations is a standard exercise, which results in the bound

$$P_\varepsilon \leq (1 + 2\theta \varepsilon)^{-\ell/2} \left( \prod\nolimits_{i \in I}(1 - 2\theta c_i) \right)^{-1/2} \quad \text{when } \theta < (2c_i)^{-1} \text{ for each } i \in I.$$

We make a coarse estimate to arrive at

$$P_\varepsilon \le (1 + 2\theta\varepsilon)^{-\ell/2} \left(1 - 2\theta \sum_{i \in I} c_i\right)^{-1/2} \quad \text{when } \theta < \left(2 \sum_{i \in I} c_i\right)^{-1}.$$

The last bound follows from repeated application of the numerical inequality $(1-s)(1-t) \ge 1 - (s+t)$, which is valid when $st \ge 0$.

Next, we must identify a suitable value for $\theta$. It is possible to minimize the probability bound with respect to $\theta$, but it is more expedient to select $\theta^{-1} = 4 \sum_{i \in I} c_i$. This choice yields

$$P_\varepsilon \le \sqrt{2} \left(1 + \frac{\varepsilon/2}{\sum_{i \in I} c_i}\right)^{-\ell/2} = \sqrt{2} \left[1 + \frac{\varepsilon/2}{\sum_{a_i < 1-\varepsilon} \varphi^2(a_i)(1 - \varepsilon - a_i)}\right]^{-\ell/2}. \tag{6.2}$$

It remains to choose a good polynomial $\varphi$. Fortunately, we have already done the work in Section 4.2.

### 6.3. **Probability Bound without a Spectral Gap.**

In this section, we use (6.2) to derive the probability bound that appears in Theorem 1.6. Set the parameter $\beta = 1 - \varepsilon$. Obviously,

$$\sum_{a_i < 1-\varepsilon} \varphi_1^2(a_i)(1 - \varepsilon - a_i) = \sum_{a_i < \beta} \varphi_1^2(a_i)(\beta - a_i).$$

To control the terms in this sum, we need second-kind Chebyshev polynomials because they satisfy the weighted uniform bound (4.8). For a partition $q = q_1 + q_2$, consider the polynomial

$$\varphi_1(s) := \frac{s^{q_1} U_{2q_2}\left(\sqrt{s/\beta}\right)}{U_{2q_2}\left(\sqrt{1/\beta}\right)} \in \mathscr{P}_q(1).$$

According to (4.10) and (4.5), this polynomial satisfies

$$\varphi_1^2(s)(\beta - s) \le \frac{4\varepsilon s^{2q_1} \delta^{2q_2+1}}{(1 - \delta^{2q_2+1})^2} \quad \text{for } 0 \le s \le \beta \quad \text{where} \quad \delta \le e^{-2\sqrt{\varepsilon}}. \tag{6.3}$$

See Section 4.2.5 for further discussion.

Using these facts, we may estimate that

$$\sum_{a_i < 1-\varepsilon} \varphi_1^2(a_i)(1 - \varepsilon - a_i) \le 4\varepsilon \left(\sum_{a_i < \beta} a_i^{2q_1}\right) \delta^{2q_2+1} \left(1 - \delta^{2q_2+1}\right)^{-2}$$

$$\le 4\varepsilon \operatorname{srk}(q_1) \delta^{2q_2+1} \left(1 - \delta^{2q_2+1}\right)^{-2}.$$

In the last step, we bound the sum in terms of the stable rank (1.11). We rely on the normalization (5.2) to recognize the stable rank.

Select $\varphi = \varphi_1$ in our probability bound (6.2). Using the last display, we arrive at

$$P_\varepsilon \le \sqrt{2} \left[1 + \frac{\left(1 - \delta^{2q_2+1}\right)^2}{8 \operatorname{srk}(q_1) \delta^{2q_2+1}}\right]^{-\ell/2}.$$

We can develop a lower bound on the bracket as follows.

$$1 + \frac{\left(1 - \delta^{2q_2+1}\right)^2}{8 \operatorname{srk}(q_1) \delta^{2q_2+1}} \ge 1 + \frac{1 - 2\delta^{2q_2+1}}{8 \operatorname{srk}(q_1) \delta^{2q_2+1}} \ge \frac{1}{8 \operatorname{srk}(q_1) \delta^{2q_2+1}}.$$

The last inequality follows from the fact that $\operatorname{srk}(q_1) \ge 1$ because $A$ is not a multiple of the identity. Combining the last two displays, we obtain

$$P_\varepsilon \le \sqrt{2} \left[8 \operatorname{srk}(q_1) \delta^{2q_2+1}\right]^{\ell/2} \le \sqrt{2} \left[8 \operatorname{srk}(q_1) \cdot e^{-2(2q_2+1)\sqrt{\varepsilon}}\right]^{\ell/2}. \tag{6.4}$$

The final relation is a consequence of the bound for $\delta$ in (6.3). This is the required statement.

6.4. **Probability Bound with a Spectral Gap.** Now, we use (6.2) to establish the probability bound that appears in Theorem 1.8. Recall that the spectral gap $\gamma$ is defined in (1.10). This time, set $\beta = 1 - \gamma$. Since $1 - \gamma$ is the first eigenvalue smaller than one and $\varepsilon > 0$,

$$\sum\nolimits_{a_i \le 1-\varepsilon} \varphi_2^2(a_i)(1-\varepsilon - a_i) \le \sum\nolimits_{a_i \le 1-\gamma} \varphi_2^2(a_i).$$

To control the sum, we need first-kind Chebyshev polynomials because they satisfy the uniform bound (4.2). For a partition $q = q_1 + q_2$, construct the polynomial

$$\varphi_2(s) := \frac{s^{q_1} T_{q_2}((2/\beta)s - 1)}{T_{q_2}((2/\beta) - 1)} \in \mathscr{P}_q(1). \tag{6.5}$$

According to (4.5) and (4.6), this polynomial satisfies

$$\varphi_2^2(s) \le 4 s^{2q_1} e^{-4q_2 \sqrt{\gamma}} \quad \text{for } 0 \le s \le \beta. \tag{6.6}$$

See Section 4.2.3 for more details.

Using these facts, we calculate that

$$\sum\nolimits_{a_i \le 1-\varepsilon} \varphi_2^2(a_i)(1-\varepsilon - a_i) \le 4 \left( \sum\nolimits_{a_i \le 1-\gamma} a_i^{2q_1} \right) e^{-4q_2 \sqrt{\gamma}} \le 4 \operatorname{srk}(q_1) \cdot e^{-4q_2 \sqrt{\gamma}}.$$

We have invoked (5.2) to identify the stable rank (1.11).

Instantiate the probability estimate (6.2) with $\varphi = \varphi_2$, and substitute in the last display to obtain

$$P_\varepsilon \le \sqrt{2} \left[ 1 + \frac{\varepsilon}{8 \operatorname{srk}(q_1) \cdot e^{-4q_2 \sqrt{\gamma}}} \right]^{-\ell/2} \le \sqrt{2} \left[ \frac{8 \operatorname{srk}(q_1)}{\varepsilon} \cdot e^{-4q_2 \sqrt{\gamma}} \right]^{\ell/2}. \tag{6.7}$$

We have used the numerical inequality $(1 + 1/s)^{-1} \le s$, valid for $s > 0$. This is the advertised result.

## 7. A BOUND FOR THE EXPECTED ERROR WITHOUT A SPECTRAL GAP

In this section, we establish the expectation bound that appears in Theorem 1.6. To obtain this result, we simply integrate the probability bound (6.4). Surprisingly, this approach appears to be more effective than a direct computation of the expected error. This insight yields a better expected error bound than the one obtained in [KW92].

7.1. **Computing the Expectation.** We may express the expectation of the relative error as an integral:

$$E := \mathbb{E} \operatorname{err}(\xi_{\max}(A; \Omega; q)) = \int_0^1 \mathbb{P} \left\{ \operatorname{err}(\xi_{\max}(A; \Omega; q)) \ge \varepsilon \right\} d\varepsilon = \int_0^1 P_\varepsilon \, d\varepsilon.$$

The limits of the integral follow from the fact that the relative error falls in the interval $[0, 1]$. We split the integral at a value $c > 0$, to be determined later. Then make the estimates

$$E \le c + \int_c^1 P_\varepsilon \, d\varepsilon \le c + \sqrt{2} \, (8 \operatorname{srk}(q_1))^{\ell/2} \int_c^\infty e^{-(2q_2+1)\ell \sqrt{\varepsilon}} \, d\varepsilon.$$

To obtain the first inequality, we use the trivial bound $P_\varepsilon \le 1$. The second inequality is a consequence of (6.4). The remaining integral can be calculated by changing the variable and integrating by parts. Indeed,

$$\int_c^\infty e^{-p\sqrt{\varepsilon}} \, d\varepsilon = 2 \left( \frac{\sqrt{c}}{p} + \frac{1}{p^2} \right) e^{-p\sqrt{c}} \quad \text{for } p > 0.$$

Together, the last two displays yield

$$E \le c + 2\sqrt{2} (8 \operatorname{srk}(q_1))^{\ell/2} \left( \frac{\sqrt{c}}{(2q_2+1)\ell} + \frac{1}{(2q_2+1)^2 \ell^2} \right) e^{-(2q_2+1)\ell \sqrt{c}}.$$

Now, select the (optimal) value

$$c = \left( \frac{\ell^{-1} \log 2 + \log(8 \operatorname{srk}(q_1))}{2(2q_2+1)} \right)^2.$$

Combine the last two displays to reach

$$E \leq \left( \frac{\ell^{-1}(2 + \log 2) + \log(8\,\mathrm{srk}(q_1))}{2(2q_2 + 1)} \right)^2$$

Bound the numerical constant by 2.70 to complete the proof.

## 8. A Bound for the Expected Error with a Spectral Gap

Last, we establish the expectation bounds for the relative error that appear in Theorem 1.8. In this case, we achieve better results by a direct computation, rather than by integrating the probability bound (6.7). These arguments are inspired by the approach in [KW92], but our computations are technically easier because the random vector $\boldsymbol{x}$ has independent entries.

### 8.1. Form of the Expected Error.
Fix a polynomial $\varphi \in \mathscr{P}_q(1)$. Take the expectation of the error bound (5.3):

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \mathbb{E}\left[ \frac{\sum_{i>1} X_i^2 \varphi^2(a_i)(1 - a_i)}{X_1^2 + \sum_{i>1} X_i^2 \varphi^2(a_i)} \right].$$

Set the parameter $\beta = 1 - \gamma$. By the definition (1.10) of the spectral gap, each eigenvalue of $\boldsymbol{A}$ that exceeds $\beta$ equals the maximum eigenvalue $a_1 = 1$. Therefore,

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \mathbb{E}\left[ \frac{\sum_{a_i \leq \beta} X_i^2 \varphi^2(a_i)}{X_1^2 + \sum_{a_i \leq \beta} X_i^2 \varphi^2(a_i)} \right]. \tag{8.1}$$

By independence, we may compute the expectation with respect to $X_1$, holding $X_i$ fixed for each $i$ where $a_i \leq \beta$. The computation of this integral depends on the block size $\ell$.

### 8.2. Error Bound for Block Size $\ell \geq 3$.
We begin with the case $\ell \geq 3$. Use the fact that $X_1 \sim \mathrm{CHI}(\ell)$ to compute that

$$\mathbb{E}\left[ \frac{1}{X_1^2 + c} \right] = \frac{1}{2} e^{c/2} \int_1^\infty s^{-\ell/2} e^{-cs/2}\,\mathrm{d}s \leq \frac{1}{(\ell - 2) + c} \quad \text{for } \ell \geq 3 \text{ and } c \geq 0.$$

The first relation depends on a standard identity for the partial gamma function [OLBC10, Sec. 8.6.4]. The second relation is a classic bound for the exponential integral due to Hopf [Hop34, p. 26]; see [Gau59] or [AS64, Sec. 5.1.19].

With $c = \sum_{a_i \leq \beta} X_i^2 \varphi^2(a_i)$, the last two displays imply that

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \mathbb{E}\left[ \frac{\sum_{a_i \leq \beta} X_i^2 \varphi^2(a_i)}{(\ell - 2) + \sum_{a_i \leq \beta} X_i^2 \varphi^2(a_i)} \right].$$

The function $t \mapsto t/(s + t)$ is concave, so Jensen's inequality allows us to draw the expectation inside the function to reach

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega}; q)) \leq \frac{\sum_{a_i \leq \beta} \varphi^2(a_i)}{(\ell - 2) + \sum_{a_i \leq \beta} \varphi^2(a_i)}.$$

Indeed, $\mathbb{E}\,X_i^2 = 1$ because $X_i$ is standard normal for each $i > 1$.

Introduce the polynomial $\varphi = \varphi_2$ from (6.5) into the last display. Using the upper bound (6.6) and the fact that $t \mapsto t/(s + t)$ is increasing, we arrive at

$$\mathbb{E}\,\mathrm{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \leq \frac{4\,\mathrm{srk}(q_1) \cdot e^{-4q_2\sqrt{\gamma}}}{(\ell - 2) + 4\,\mathrm{srk}(q_1) \cdot e^{-4q_2\sqrt{\gamma}}}.$$

This is the advertised result for $\ell = 3$ in Theorem 1.8.

8.3. **Error Bound for Block Size** $\ell = 2$. Now, assume that the block size $\ell = 2$. In this case, $X_1 \sim \text{CHI}(2)$. Therefore, it holds that

$$\mathbb{E}\left[\frac{1}{X_1^2 + c}\right] = \frac{1}{2}e^{c/2}\int_1^\infty s^{-1}e^{-cs/2}\,\mathrm{d}s \le \frac{1}{2}\log\left(1 + \frac{2}{c}\right) \quad \text{for } c > 0.$$

The first relation is an immediate consequence of the definition of the chi-square density and a change of variables. The second relation is a classic bound for the exponential integral [AS64, Sec. 5.1.20].

The rest of the argument follows the same path as in the case $\ell \ge 3$. With $c = \sum_{a_i \le \beta} X_i^2 \varphi^2(a_i)$, we combine (8.1) with the last display to obtain

$$\mathbb{E}\,\text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \le \mathbb{E}\left[\frac{\sum_{a_i \le \beta} X_i^2 \varphi^2(a_i)}{2}\cdot\log\left(1 + \frac{2}{\sum_{a_i \le \beta} X_i^2 \varphi^2(a_i)}\right)\right].$$

The function $t \mapsto t\log(1 + 1/t)$ is concave and increasing, so Jensen's inequality yields

$$\mathbb{E}\,\text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \le \frac{\sum_{a_i \le \beta} \varphi^2(a_i)}{2}\cdot\log\left(1 + \frac{2}{\sum_{a_i \le \beta} \varphi^2(a_i)}\right).$$

Select the polynomial $\varphi = \varphi_2$ from (6.5), and invoke the bound (6.6) to arrive at

$$\mathbb{E}\,\text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \le \frac{1}{2}\cdot 4\,\text{srk}(q_1)e^{-4q_2\sqrt{\gamma}}\cdot\log\left(1 + \frac{2}{4\,\text{srk}(q_1)e^{-4q_2\sqrt{\gamma}}}\right).$$

This is the desired outcome for block size $\ell = 2$ in Theorem 1.8.

8.4. **Error Bound for Block Size** $\ell = 1$. Finally, we complete the analysis for the case where the block size $\ell = 1$. Use the fact that $X_1 \sim \text{CHI}(1)$ to see that

$$\mathbb{E}\left[\frac{1}{X_1^2 + c}\right] = \frac{1}{\sqrt{c}}e^{c/2}\int_{\sqrt{c}}^\infty e^{-s^2/2}\,\mathrm{d}s \le \sqrt{\frac{2\pi}{c}}.$$

The first relation follows from [OLBC10, Sec. 8.6.4] after a change of variable. The second relation is a well-known tail bound for a standard normal random variable. Continuing as we have done before,

$$\mathbb{E}\,\text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \le \sqrt{2\pi}\,\mathbb{E}\left[\sum\nolimits_{a_i \le \beta} X_i^2 \varphi^2(a_i)\right]^{1/2} \le \sqrt{2\pi}\left[\sum\nolimits_{a_i \le \beta} \varphi^2(a_i)\right]^{1/2}.$$

Insert the bound (6.6) for the polynomial $\varphi = \varphi_2$ to arrive at

$$\mathbb{E}\,\text{err}(\xi_{\max}(\boldsymbol{A}; \boldsymbol{\Omega})) \le \sqrt{8\pi\,\text{srk}(q_1)}\cdot e^{-2q_2\sqrt{\gamma}}.$$

Of course, the relative error is always bounded above by one. This is the result stated in Theorem 1.8.

## REFERENCES

[AS64]     M. Abramowitz and I. A. Stegun. *Handbook of mathematical functions with formulas, graphs, and mathematical tables*, volume 55 of *National Bureau of Standards Applied Mathematics Series*. U.S. Government Printing Office, Washington, D.C., 1964.

[BDD+00]   Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst, editors. *Templates for the solution of algebraic eigenvalue problems: A practical guide*, volume 11 of *Software, Environments, and Tools*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2000.

[BES05]    C. A. Beattie, M. Embree, and D. C. Sorensen. Convergence of polynomial restart Krylov methods for eigenvalue computations. *SIAM Rev.*, 47(3):492–515, 2005.

[Bha97]    R. Bhatia. *Matrix analysis*, volume 169 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, 1997.

[BT87]    J. Bourgain and L. Tzafriri. Invertibility of "large" submatrices with applications to the geometry of Banach spaces and harmonic analysis. *Israel J. Math.*, 57(2):137–224, 1987.

[BT91]    J. Bourgain and L. Tzafriri. On a problem of Kadison and Singer. *J. Reine Angew. Math.*, 420:1–43, 1991.

[CD74]    J. Cullum and W. E. Donath. A block generalization of the $s$-step Lanczos algorithm. Report RC 4845 (21570), IBM Thomas J. Watson Research Center, Yorktown Heights, New York, 1974.

[Cho87]   A. W. Chou. On the optimality of Krylov information. *J. Complexity*, 3(1):26–40, 1987.

[DCM97]   G. M. Del Corso and G. Manzini. On the randomized error of polynomial methods for eigenvector and eigenvalue estimates. *J. Complexity*, 13(4):419–456, 1997.

[DI19]    P. Drineas and I. C. F. Ipsen. Low-rank matrix approximations do not need a singular value gap. *SIAM J. Matrix Anal. Appl.*, 40(1):299–319, 2019.

[DIKMI18] P. Drineas, I. C. F. Ipsen, E.-M. Kontopoulou, and M. Magdon-Ismail. Structural convergence results for approximation of dominant subspaces from block Krylov spaces. *SIAM J. Matrix Anal. Appl.*, 39(2):567–586, 2018.

[Dix83]   J. D. Dixon. Estimating extremal eigenvalues and condition numbers of matrices. *SIAM J. Numer. Anal.*, 20(4):812–814, 1983.

[DKM06]   P. Drineas, R. Kannan, and M. W. Mahoney. Fast Monte Carlo algorithms for matrices. II. Computing a low-rank approximation to a matrix. *SIAM J. Comput.*, 36(1):158–183, 2006.

[FKV98]   A. Frieze, R. Kannan, and S. Vempala. Fast Monte Carlo algorithms for finding low-rank approximations. In *Proc. 39th Ann. IEEE Symp. Foundations of Computer Science (FOCS)*, pages 370–378, 1998.

[Gau59]   W. Gautschi. Exponential integral $\int_1^\infty e^{-xt} t^{-n} dt$ for large values of $n$. *J. Res. Nat. Bur. Standards*, 62:123–125, 1959.

[GU77]    G. H. Golub and R. Underwood. The block Lanczos method for computing eigenvalues. In *Mathematical software, III (Proc. Sympos., Math. Res. Center, Univ. Wisconsin, Madison, Wis., 1977)*, pages 361–377. Publ. Math. Res. Center, No. 39. Academic Press, New York, 1977.

[Gu15]    M. Gu. Subspace iteration randomization and singular value problems. *SIAM J. Sci. Comput.*, 37(3):A1139–A1173, 2015.

[GVL13]   G. H. Golub and C. F. Van Loan. *Matrix computations*. Johns Hopkins Studies in the Mathematical Sciences. Johns Hopkins University Press, Baltimore, MD, fourth edition, 2013.

[HMST11]  N. Halko, P.-G. Martinsson, Y. Shkolnisky, and M. Tygert. An algorithm for the principal component analysis of large data sets. *SIAM J. Sci. Comput.*, 33(5):2580–2594, 2011.

[HMT11]   N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.

[Hop34]   E. Hopf. *Mathematical problems of radiative equilibrium*. Number 31 in Cambridge Tracts in Mathematics and Mathematical Physics. Cambridge Univ. Press, 1934.

[KW92]    J. Kuczyński and H. Woźniakowski. Estimating the largest eigenvalue by the power and Lanczos algorithms with a random start. *SIAM J. Matrix Anal. Appl.*, 13(4):1094–1122, 1992.

[KW94]    J. Kuczyński and H. Woźniakowski. Probabilistic bounds on the extremal eigenvalues and condition number by the Lanczos algorithm. *SIAM J. Matrix Anal. Appl.*, 15(2):672–691, 1994.

[Lan50]   C. Lanczos. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *J. Research Nat. Bur. Standards*, 45:255–282, 1950.

[LS05]    J. Liesen and Z. Strakoš. GMRES convergence analysis for a convection-diffusion model problem. *SIAM J. Sci. Comput.*, 26(6):1989–2009, 2005.

[LW98]    Z. Leyk and H. Woźniakowski. Estimating a largest eigenvector by Lanczos and polynomial algorithms with a random start. *Numer. Linear Algebra Appl.*, 5(3):147–164, 1998.

[LZ15]    R.-C. Li and L.-H. Zhang. Convergence of the block Lanczos method for eigenvalue clusters. *Numer. Math.*, 131(1):83–113, 2015.

[Meu06]   G. Meurant. *The Lanczos and conjugate gradient algorithms: From theory to finite precision computations*, volume 19 of *Software, Environments, and Tools*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2006.

[MM15]    C. Musco and C. Musco. Stronger and faster approximate singular value decomposition via the block Lanczos method. In *Proc. Adv. Neural Information Processing Systems 28 (NIPS 2015)*, 2015. Available at http://arXiv.org/abs/1504.05477.

[MRT11]   P.-G. Martinsson, V. Rokhlin, and M. Tygert. A randomized algorithm for the decomposition of matrices. *Appl. Comput. Harmon. Anal.*, 30(1):47–68, 2011.

[MT20]    P.-G. Martinsson and J. A. Tropp. Randomized numerical linear algebra: Foundations and algorithms. *Acta Numerica*, pages 403–472, 2020.

[Nem91]   A. S. Nemirovsky. On optimality of Krylov's information when solving linear operator equations. *J. Complexity*, 7(2):121–130, 1991.

[NY83]    A. S. Nemirovsky and D. B. Yudin. *Problem complexity and method efficiency in optimization*. A Wiley-Interscience Publication. John Wiley & Sons, Inc., New York, 1983. Translated from the Russian and with a preface by E. R. Dawson, Wiley-Interscience Series in Discrete Mathematics.

[OLBC10]  F. W. J. Olver, D. W. Lozier, R. F. Boisvert, and C. W. Clark, editors. *NIST handbook of mathematical functions*. U.S. Department of Commerce, National Institute of Standards and Technology, Washington, DC; Cambridge University Press, Cambridge, 2010.

[Pai71]  C. C. Paige. *The computation of eigenvalues and eigenvectors of very large sparse matrices*. PhD thesis, Univ. London, 1971.

[Par98]  B. N. Parlett. *The symmetric eigenvalue problem*, volume 20 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1998. Corrected reprint of the 1980 original.

[PRTV98]  C. H. Papadimitriou, P. Raghavan, H. Tamaki, and S. Vempala. Latent semantic indexing: A probabilistic analysis. In *Proc. 17th ACM Symp. Principles of Database Systems (PODS)*, 1998.

[RST09]  V. Rokhlin, A. Szlam, and M. Tygert. A randomized algorithm for principal component analysis. *SIAM J. Matrix Anal. Appl.*, 31(3):1100–1124, 2009.

[Saa80]  Y. Saad. On the rates of convergence of the Lanczos and the block-Lanczos methods. *SIAM J. Numer. Anal.*, 17(5):687–706, 1980.

[Saa11]  Y. Saad. *Numerical methods for large eigenvalue problems*, volume 66 of *Classics in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2011. Revised edition of the 1992 original.

[SAR17]  M. Simchowitz, A. E. Alaoui, and B. Recht. On the gap between strict-saddles and true convexity: An $\omega(\log d)$ lower bound for eigenvector approximation. Available at http://arXiv.org/abs/1704.04548, Apr. 2017.

[SEAR18]  M. Simchowitz, A. El Alaoui, and B. Recht. Tight query complexity lower bounds for PCA via finite sample deformed Wigner law. In *STOC'18—Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing*, pages 1249–1259. ACM, New York, 2018.

[Tro18]  J. A. Tropp. Analysis of randomized block Krylov methods. ACM Report 2018-02, California Institute of Technology, 2018.

[TW21]  J. A. Tropp and R. J. Webber. Randomized algorithms for low-rank matrix approximation: Design, analysis, and applications. In preparation, Oct. 2021.

[vdSvdV86]  A. van der Sluis and H. A. van der Vorst. The rate of convergence of conjugate gradients. *Numer. Math.*, 48(5):543–560, 1986.

[WZZ15]  S. Wang, Z. Zhang, and T. Zhang. Improved analyses of the randomized power method and block Lanczos method. Available at http://arXiv.org/abs/1508.06429, Aug. 2015.

[YGL18]  Q. Yuan, M. Gu, and B. Li. Superlinear convergence of randomized block Lanczos algorithm. In *2018 IEEE International Conference on Data Mining (ICDM)*, pages 1404–1409, 2018.