

# PROBABILISTIC ALGORITHMS FOR FINDING MATRIX DECOMPOSITIONS

ALEX NOWAK

**ABSTRACT.** Low-rank matrix approximations are oblique in many areas ranging from data analysis to scientific computing. From a data science point of view, probably the most important application is due to Principal Component Analysis (PCA), which aims to reveal hidden linear structure in massive datasets through a low-rank matrix decomposition. Consequently, the complexity of the algorithm plays a central role in the applicability of the algorithms to big data. The most common approximative factorization is the so-called truncated singular value decomposition (k-SVD) which can be computed in  $\mathcal{O}(mnk)$  floating-point operations, where  $k$  is the target rank of the decomposition and  $m$  and  $n$  are the corresponding dimensions of the matrix. In this review, we introduce to the reader randomized algorithms that can achieve the aforementioned task with numerous advantages compared to the classical algorithms. These algorithms are based on the fact that the image of a low-rank matrix can be approximated by the action of the matrix to a reasonable amount of random vectors from the input space. Starting from this point, it is possible to develop algorithms that achieve a complexity of  $\mathcal{O}(mn \log k)$  for dense-matrices, matches the flop count of classical Krylov subspace methods for sparse matrices with a gain in robustness, and for large matrices that can not be stored in memory (RAM), they achieve a constant number of passes compared to the  $\mathcal{O}(k)$  for classical algorithms.

## INTRODUCTION

Matrix factorization is listed as one of the most influential set of techniques during the 20th century [3], among the Fast Fourier Transform, MCMC sampling methods and others. As Stewart [7] argues, the principle of the decompositional approach aims to construct computational platforms from which a variety of problems can be solved.

Although the decompositional approach to matrix computation remains fundamental, nowadays in the era of big data, most of the classical algorithms are inadequate to tackle most of the problems.

The empirical covariance matrices derived from datasets are now incredibly big, making most of the classical approaches too expensive in terms of computation. Moreover, it is common in information sciences to have data which is missing or inaccurate. This gives the opportunity to sacrifice some accuracy on the algorithm to gain on computation, which classical algorithms are not able to do. Another important aspect is the role of data transfer in the computational cost of a given algorithm, i.e., techniques that may perform fewer passes over the data may be substantially faster in practice.

In this review, we present, analyse and test *randomized* algorithms for matrix factorizations. This set of novel techniques addresses the issues stated above, i.e., can trade-off computation and accuracy to an arbitrary precision, gain on robustness, and reduce the number of passes on big datasets when data transfer is expensive.

The purpose of approximated low-rank matrix factorization is to factorize a given matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$  into a product of two smaller matrices  $\mathbf{B} \in \mathbb{R}^{m \times k}$  and  $\mathbf{C} \in \mathbb{R}^{k \times n}$ .

$$(0.1) \quad \begin{matrix} \mathbf{A} \\ m \times n \end{matrix} \approx \begin{matrix} \mathbf{B} \\ m \times k \end{matrix} \begin{matrix} \mathbf{C} \\ k \times n \end{matrix},$$

The matrix  $\mathbf{B} \times \mathbf{C}$  in 0.1 is called a rank- $k$  approximation of the matrix  $\mathbf{A}$ .

The inner dimension  $k$  is called the *numerical rank* of the matrix. This quantity differs from the *algebraic rank*, which is defined as the dimension of the image. The numerical rank is commonly defined as follows

$$(0.2) \quad r(\mathbf{A}) := \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|^2} = \sum_{j=1}^{\min(m,n)} \left( \frac{\sigma_j}{\sigma_1} \right)^2$$

and it gives a better understanding of how accurate a rank- $k$  approximation can be. Note that we always have  $r(\mathbf{A}) \leq \text{rank}(\mathbf{A})$ . The notion of numerical rank appears in [8] and has been studied at the Theory Reading Group course of the master in depth <sup>1</sup>.

The task of computing a low-rank approximation to a given matrix can be split into two computational stages. The first is to construct a low dimensional subspace that can capture the action of the matrix. The second is to restrict the matrix to the low dimensional subspace and then compute a standard factorization (QR, SVD, etc) of the reduced matrix.

- **Stage 1:** Compute an approximate basis for the range of the input matrix  $\mathbf{A}$ . We want to find a matrix  $\mathbf{Q}$  with a small number of orthonormal columns such that

$$(0.3) \quad \mathbf{A} \approx \mathbf{Q}\mathbf{Q}^* \mathbf{A}$$

The main idea is to approximate the range of the matrix via a randomized method. This can be accomplished by iteratively computing the image of random vectors from the input space and then orthogonalizing. All the randomness of the algorithms will belong to this first stage.

- **Stage 2:** Given a matrix  $\mathbf{Q}$  that satisfies 0.3, compute a standard factorization (QR, SVD, etc.) of  $\mathbf{A}$ . Note that taking  $\mathbf{B} = \mathbf{Q}$  and  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$  we already have a low rank approximation of the matrix. There is no randomness at this stage and only classical linear algebra computations are involved.

During the rest of the introduction, we will review some basics about matrix approximation and we will provide to the reader the basic aspects and insights of both stages, which will be further studied in depth in the main body.

**0.1. Approximating the range of a matrix via randomness.** The problem of finding the best  $\epsilon$ -approximation of a given matrix  $\mathbf{A}$  is called the *fixed-precision problem*. More concretely, we are given a tolerance  $\epsilon$  and the goal is to find a matrix  $\mathbf{Q}$  with  $k = k(\epsilon)$  columns such that

$$(0.4) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \epsilon$$

The goal here is to find a  $\mathbf{Q}$  with the smaller number of columns possible.

Another closely related problem is the so-called *fixed-rank problem*, which seeks to find the best rank- $k$  approximation of the matrix.

$$(0.5) \quad \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\|.$$

The Singular Value Decomposition (SVD) is key to analyze this problem. Recall that the SVD of a matrix  $\mathbf{A}$  is the following decomposition

$$(0.6) \quad \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \sum_{j=1}^{\text{rank}(\mathbf{A})} \sigma_j \mathbf{u}_j \mathbf{v}_j^*$$

where  $\{\mathbf{u}_k\}_k, \{\mathbf{v}_k\}_k$  are orthonormal basis on the output and input space respectively,  $\sigma_1 \geq \sigma_2, \dots, \sigma_{\text{rank}(\mathbf{A})} \geq 0$  are the ordered singular values.

<sup>1</sup>The comparison between both notions of rank can be better understood through the following characterization.  $\text{rank}(\mathbf{A}) = \dim(\mathbf{A}B_2^n)$  and  $r(\mathbf{A}) = d(\mathbf{A}B_2^n)$  where  $B_2^n$  is the euclidean ball, hence,  $\mathbf{A}B_2^n$  is the ellipsoid with the axis of magnitude the singular values  $\sigma_j$ 's. Here,  $\dim(\cdot)$  denotes the algebraic dimension and  $d(\cdot)$  denotes the *statistical dimension*. The statistical dimension is defined as  $d(T) := \frac{h(T-T)^2}{\text{diam}(T)^2} \sim \frac{w(T)^2}{\text{diam}(T)^2}$  where  $h(T)^2 = \mathbb{E} \sup_{t \in T} \langle g, t \rangle^2$  and  $w(T) = \mathbb{E} \sup_{t \in T} \langle g, t \rangle$  is the gaussian width. As discussed in [8], the statistical dimension is a more stable notion of dimension, in the same way that the numerical rank is more stable than the algebraic rank.

The SVD provides an optimal answer to the *fixed precision problem* [6] through the following important observation <sup>2</sup>.

$$(0.7) \quad \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\| = \sigma_{j+1}.$$

It is straightforward to check that the optimum is attained at  $\mathbf{X}_* = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^*$ , namely, the  $k$ -truncated SVD of the matrix  $\mathbf{A}$ . More precisely, we have that  $\mathbf{B} = \mathbf{U} \Sigma_{[j]}^{1/2}$  and  $\mathbf{C} = \Sigma_{[j]}^{1/2} \mathbf{V}^*$  are the best solutions 0.1 when the rank is fixed.

Let's suppose now that we know the desired rank  $k$  in advance. The goal is to find a matrix  $\mathbf{Q}$  with  $k + p$  orthonormal columns such that

$$(0.8) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \approx \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\|$$

where  $p$  is called the oversampling parameter.

EXPLICACIO

**0.2. Intuition of the randomized method to find  $\mathbf{Q}$ .** The key observation that leverage these methods is the fact that the matrix  $\mathbf{Q}$  can be found by sampling, and now the reader will obtain intuition on how this can be done.

Suppose we seek a basis for the range of a matrix  $\mathbf{A}$  with algebraic rank  $k$ . Random elements  $\mathbf{y}^{(i)}$  from the range can be computed by computing the image of random vectors  $\mathbf{w}^{(i)}$  of the input space. Let us repeat this process  $k$  times:

$$(0.9) \quad \mathbf{y}^{(i)} = \mathbf{A} \mathbf{w}^{(i)}, \quad i = 1, \dots, k$$

Thanks to the randomness, the set  $\{\mathbf{w}^{(i)}\}_{i=1}^k$  is likely to be in general linear position and no vector will fall in  $\ker \mathbf{A}$  if this is a set of measure zero under the probability measure we sampled from. Therefore, an orthogonalization procedure gives the desired orthonormal basis.

What happens if the matrix  $\mathbf{A}$  has not exact algebraic rank equal to  $k$ ? Write  $\mathbf{A} = \mathbf{B} + \mathbf{E}$  where  $\mathbf{B}$  is a rank- $k$  matrix containing the information we seek and  $\mathbf{E}$  a small perturbation. We want a basis that covers the range of  $\mathbf{B}$ , however, if we repeat procedure 0.10, the vectors will be affected by the perturbation and the  $\{\mathbf{y}^{(i)}\}_{i=1}^k$  will have small components that will make them fall outside the desired space.

To overcome this issue, the idea is to take  $p$  more samples:

$$(0.10) \quad \mathbf{y}^{(i)} = \mathbf{A} \mathbf{w}^{(i)} = \mathbf{B} \mathbf{w}^{(i)} + \mathbf{E} \mathbf{w}^{(i)}, \quad i = 1, \dots, k + p$$

The enriched set  $\{\mathbf{y}^{(i)}\}_{i=1}^{k+p}$  has much more chance of spanning the desired subspace, and this is grounded with some theoretical results holding in high probability. The theory also shows that  $p$  can be quite small. In practice,  $p = 5$  is more than enough.

PROTO-ALGORITHM: SOLVING THE FIXED-RANK PROBLEM

- 1 Draw a random  $n \times (k + p)$  test matrix  $\mathbf{\Omega}$ .
- 2 Form the matrix product  $\mathbf{Y} = \mathbf{A} \mathbf{\Omega}$ .
- 3 Construct a matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ .

<sup>2</sup>This is a direct consequence of the Courant-Fisher's *min-max theorem*. This theorem offers a variational characterization of the singular values of a given matrix  $\mathbf{A}$ .

$$\sigma_i(\mathbf{A}) = \min_{E, \dim E = i} \max_{\mathbf{x} \in S(E)} \|\mathbf{A} \mathbf{x}\|$$

**0.3. Construction of standard matrix factorizations from  $\mathbf{Q}$ .** This corresponds to **Stage 2** of the algorithm. Once we have  $\mathbf{Q}$  such that  $\mathbf{A} \approx \mathbf{Q}\mathbf{Q}^*\mathbf{A}$ , taking  $\mathbf{B} = \mathbf{Q}$  and  $\mathbf{C} = \mathbf{Q}^*\mathbf{A}$  we produce a low rank matrix decomposition  $\mathbf{A} \approx \mathbf{BC}$ .

A lot of questions need to be addressed in order to turn these methods into a technology or *off-the-shelf* algorithms, namely, theoretical guarantees on the accuracy of the approximate factorization and experimentally study the gap between theory and practice, i.e, how these methods perform in practice.

## ALGORITHMS

In this section we will describe the randomized algorithms in detail, provide the corresponding computational complexity analysis, and state the main theoretical results that guarantee the accuracy of the approximation.

### 1. STAGE 1

In the introductory section, we provided some intuition on the randomized procedure and we developed a general Proto-Algorithm 0.2 to find the matrix  $\mathbf{Q}$ . However, Proto-Algorithm 0.2 is very general and can be tuned depending on the problem requirements. The number  $T_{\text{basic}}$  of flops required by Proto-Algorithm 0.2 satisfies

$$(1.1) \quad T_{\text{basic}} \sim \ell n T_{\text{rand}} + \ell T_{\text{mult}} + \ell^2 m$$

where  $T_{\text{rand}}$  is the cost of generating a Gaussian random number and  $T_{\text{mult}}$  is the cost of multiplying  $\mathbf{A}$  by a vector.

We will now describe some specific realizations of Proto-Algorithm 0.2 that will be intended for problems with different requirements.

**1.0.1. Randomized Range Finder.** This is the most naive and simplest implementation of Proto-Algorithm 0.2. Given an oversampling parameter  $p$ , the *Randomized Range Finder* performs Proto-Algorithm 0.2 with a gaussian test matrix  $\mathbf{\Omega} \in \mathbb{R}^{n \times \ell}$  with  $\ell = k + p$  and  $k$  being a pre-specified target rank. Then, it orthogonalizes the rows of the resulting matrix  $\mathbf{Y}$  by computing a QR decomposition. A numerical issue arises when computing the orthogonalization procedure due to the fact that the columns of  $\mathbf{Y}$  are almost linearly dependent. The authors in [5] found that using the *double orthogonalization* [1] was enough to guarantee stability of the procedure.

RANDOMIZED RANGE FINDER	
1	Draw an $n \times \ell$ Gaussian random matrix $\mathbf{\Omega}$ .
2	Form the $m \times \ell$ matrix $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ .
3	Construct an $m \times \ell$ matrix $\mathbf{Q}$ whose columns form an orthonormal basis for the range of $\mathbf{Y}$ , e.g., using the QR factorization $\mathbf{Y} = \mathbf{QR}$ .

The complexity analysis of the Algorithm 1.0.1 gives

$$(1.2) \quad T_{\text{Randomized Range Finder}} \sim \mathcal{O}(mn\ell)$$

this is because generating a gaussian random number is  $\mathcal{O}(1)$  and computing a matrix vector multiplication is  $\mathcal{O}(mn)$ .

**1.0.2. Adaptive Randomized Range Finder.** One important pitfall of the *Randomized Range Finder* 1.0.1 is that it requires to know in advance the target rank  $k$ . However, if we intend to solve the *fixed-precision problem*, we need a scheme to estimate the error  $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|$  during the algorithm in order to match the required tolerance  $\epsilon$ .

This scheme is possible and it is direct consequence of the following lemma.

**Lemma 1.1.** *Let  $\mathbf{B}$  be a real  $m \times n$  matrix. Fix a positive integer  $r$  and a real number  $\alpha > 1$ . Draw an independent family  $\{\boldsymbol{\omega}^{(i)} : i = 1, 2, \dots, r\}$  of standard Gaussian vectors. Then*

$$\|\mathbf{B}\| \leq \alpha \sqrt{\frac{2}{\pi}} \max_{i=1, \dots, r} \|\mathbf{B}\boldsymbol{\omega}^{(i)}\|$$

except with probability  $\alpha^{-r}$ .

Lemma 1.1 says that we can bound the error with high probability using inexpensive computations in an online manner. The Lemma 1.1 applied to our problem reads

$$(1.3) \quad \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\| \leq 10\sqrt{\frac{2}{\pi}} \max_{i=1, \dots, r} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\boldsymbol{\omega}^{(i)}\|$$

with probability at least  $1 - 10^{-r}$ .

The high probability bound 1.3 gives a simple online scheme to decide when we have a good enough  $\mathbf{Q}$  that matches the pre-specified tolerance. The goal here is to find an integer  $l$  and a  $m \times l$  orthonormal matrix  $\mathbf{Q}^{(l)}$  such that

$$(1.4) \quad \|(\mathbf{I} - \mathbf{Q}^{(l)}(\mathbf{Q}^{(l)})^*)\mathbf{A}\| \leq \varepsilon.$$

We call *Adaptive Randomized Range Finder* 1.0.2 to the algorithm derived from Lemma 1.1 that solves this problem.

ADAPTIVE RANDOMIZED RANGE FINDER

- 1 Draw standard Gaussian vectors  $\boldsymbol{\omega}^{(1)}, \dots, \boldsymbol{\omega}^{(r)}$  of length  $n$ .
- 2 For  $i = 1, 2, \dots, r$ , compute  $\mathbf{y}^{(i)} = \mathbf{A}\boldsymbol{\omega}^{(i)}$ .
- 3  $j = 0$ .
- 4  $\mathbf{Q}^{(0)} = [\ ]$ , the  $m \times 0$  empty matrix.
- 5 **while**  $\max \left\{ \|\mathbf{y}^{(j+1)}\|, \|\mathbf{y}^{(j+2)}\|, \dots, \|\mathbf{y}^{(j+r)}\| \right\} > \varepsilon / (10\sqrt{2/\pi})$ ,
- 6      $j = j + 1$ .
- 7     Overwrite  $\mathbf{y}^{(j)}$  by  $(\mathbf{I} - \mathbf{Q}^{(j-1)}(\mathbf{Q}^{(j-1)})^*)\mathbf{y}^{(j)}$ .
- 8      $\mathbf{q}^{(j)} = \mathbf{y}^{(j)} / \|\mathbf{y}^{(j)}\|$ .
- 9      $\mathbf{Q}^{(j)} = [\mathbf{Q}^{(j-1)} \ \mathbf{q}^{(j)}]$ .
- 10    Draw a standard Gaussian vector  $\boldsymbol{\omega}^{(j+r)}$  of length  $n$ .
- 11     $\mathbf{y}^{(j+r)} = (\mathbf{I} - \mathbf{Q}^{(j)}(\mathbf{Q}^{(j)})^*)\mathbf{A}\boldsymbol{\omega}^{(j+r)}$ .
- 12    **for**  $i = (j+1), (j+2), \dots, (j+r-1)$ ,
- 13       Overwrite  $\mathbf{y}^{(i)}$  by  $\mathbf{y}^{(i)} - \mathbf{q}^{(j)} \langle \mathbf{q}^{(j)}, \mathbf{y}^{(i)} \rangle$ .
- 14    **end for**
- 15 **end while**
- 16  $\mathbf{Q} = \mathbf{Q}^{(j)}$ .

One important question regarding Algorithm 1.0.2 is how good the bound given by Lemma 1.1 is in practice. If there is a significant gap between theory and practice the optimal  $l$  will be overestimated. This question will be addressed in the experimental section 3.3.4.

**1.0.3. Randomized Power Iteration.** The *Randomized Range Finder* 1.0.1 algorithm assumes that the singular values of the matrix decay fast. This can be seen from equation 0.10, where the small singular values interfere with the calculation of the range. This intuition is made precise in Theorem REF, where the error of the approximation depends on the  $k + 1$ -th singular value.

The goal here is to reduce the weight of the small singular values by taking powers of the matrix whose range we want to approximate. Instead of applying the sampling scheme to  $\mathbf{A}$ , we will apply it to  $\mathbf{B} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}$  where  $q > 0$  is a small integer.

The matrix  $\mathbf{B}$  has the same singular vectors than  $\mathbf{A}$  (hence, the same range), but its singular values decay much more quickly.

$$(1.5) \quad \sigma_j(\mathbf{B}) = \sigma_j(\mathbf{A})^{2q+1}, \quad j = 1, 2, 3, \dots$$

The *Randomized Power Iteration* 1.0.3 algorithm is the same as the *Randomized Range Finder* 1.0.1 but replacing the formula  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  by  $\mathbf{Y} = \mathbf{B}\mathbf{\Omega}$ .

#### RANDOMIZED POWER ITERATION

- 1 Draw an  $n \times \ell$  Gaussian random matrix  $\mathbf{\Omega}$ .
- 2 Form the  $m \times \ell$  matrix  $\mathbf{Y} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}\mathbf{\Omega}$  via alternating application of  $\mathbf{A}$  and  $\mathbf{A}^*$ .
- 3 Construct an  $m \times \ell$  matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ , e.g., via the QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

The computational complexity of the algorithm is essentially the same because it only requires  $2q+1$  as many matrix-multiplications as Algorithm 1.0.1 but the number  $q$  is in practice 2,3 or 4. This can be seen from THEOREM, which shows that the power iteration drives the approximation gap to 1 exponentially fast.

1.0.4. *Fast Randomized Range Finder.* A simple inspection to equation 1.1 reveals the computational bottleneck of the sampling procedure. This is the matrix multiplication  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  that takes  $\mathcal{O}(mn\ell)$  operations for dense matrices, which is the same as the  $\ell$ -SVD (computed after a prior rank-revealing QR factorization [4]).

The key idea is to use a *structured* random matrix that allows us to compute the product in  $\mathcal{O}(mn \log(\ell))$  operations.

The simplest structured random matrix that meets our goals is the so-called *subsampled random Fourier transform* (SRFT).

An SRFT is an  $n \times \ell$  matrix of the form

$$(1.6) \quad \mathbf{\Omega} = \sqrt{\frac{n}{\ell}} \mathbf{D}\mathbf{F}\mathbf{R},$$

where

- $\mathbf{D}$  is an  $n \times n$  diagonal matrix whose entries are independent random variables uniformly distributed on the complex unit circle,
- $\mathbf{F}$  is the  $n \times n$  unitary discrete Fourier transform (DFT), whose entries take the values  $f_{pq} = n^{-1/2} e^{-2\pi i(p-1)(q-1)/n}$  for  $p, q = 1, 2, \dots, n$ , and
- $\mathbf{R}$  is an  $n \times \ell$  matrix that samples  $\ell$  coordinates from  $n$  uniformly at random, i.e., its  $\ell$  columns are drawn randomly without replacement from the columns of the  $n \times n$  identity matrix.

Now, via a subsampled FFT [9], we can compute the sample matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  with  $\mathcal{O}(mn \log(\ell))$  operations.

The total number of operations required by this procedure is reduced to

$$(1.7) \quad T_{\text{struct}} \sim mn \log(\ell) + \ell^2 n$$

Hence, the computational complexity of the approach is essentially  $\mathcal{O}(mn \log(\ell))$ .

#### FAST RANDOMIZED RANGE FINDER

- 1 Draw an  $n \times \ell$  SRFT test matrix  $\mathbf{\Omega}$ , as defined by (1.6).
- 2 Form the  $m \times \ell$  matrix  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  using a (subsampled) FFT.
- 3 Construct an  $m \times \ell$  matrix  $\mathbf{Q}$  whose columns form an orthonormal basis for the range of  $\mathbf{Y}$ , e.g., using the QR factorization  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

## 2. STAGE 2

The output of the Stage 1 produces an orthonormal matrix  $\mathbf{Q}$  whose range captures the action of the matrix  $\mathbf{A}$ . The goal of Stage 2 is to produce standard approximate matrix factorizations of  $\mathbf{A}$  using this  $\mathbf{Q}$ .

This subsection is divided into three parts; first, we will show how to compute standard approximate matrix factorizations (SVD and QR) from a general approximate low-rank factorization. Recall that taking  $\mathbf{B} = \mathbf{Q}$  and  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$  we readily have a factorization that satisfies  $\|\mathbf{A} - \mathbf{BC}\| \leq \varepsilon$ . Then, we will describe in detail the *Direct SVD* algorithm, which will consist in constructing an SVD from  $\mathbf{B}$  and  $\mathbf{C}$ . Finally, we will comment on other more involved methods that avoid computing the expensive product  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$ .

2.0.1. *Compute standard QR and SVD from a general factorization.* Now we will specify how we can compute the standards SVD and QR decompositions from a general low rank decomposition  $\|\mathbf{A} - \mathbf{BC}\| \leq \varepsilon$  maintaining the tolerance  $\varepsilon$  from Stage 1.

- *SVD decomposition:*  $\|\mathbf{A} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*\| \leq \varepsilon$ 
  - (1) Compute a QR factorization of  $\mathbf{B}$  so that  $\mathbf{B} = \mathbf{Q}_1 \mathbf{R}_1$ .
  - (2) Form the product  $\mathbf{D} = \mathbf{R}_1 \mathbf{C}$ , and compute an SVD:  $\mathbf{D} = \mathbf{U}_2 \mathbf{\Sigma} \mathbf{V}^*$ .
  - (3) Form the product  $\mathbf{U} = \mathbf{Q}_1 \mathbf{U}_2$ .
- *QR decomposition:*  $\|\mathbf{A} - \mathbf{QR}\| \leq \varepsilon$ 
  - (1) Compute a QR factorization of  $\mathbf{B}$  so that  $\mathbf{B} = \mathbf{Q}_1 \mathbf{R}_1$ .
  - (2) Form the product  $\mathbf{D} = \mathbf{R}_1 \mathbf{C}$ , and compute a QR factorization:  $\mathbf{D} = \mathbf{Q}_2 \mathbf{R}$ .
  - (3) Form the product  $\mathbf{Q} = \mathbf{Q}_1 \mathbf{Q}_2$ .

2.0.2. *Direct SVD.* The procedure described in 2.0.1 to compute the approximate SVD decomposition without sacrificing error, defines what we call the *Direct SVD* 2.0.2.

DIRECT SVD	
1	Form the matrix $\mathbf{B} = \mathbf{Q}^* \mathbf{A}$ .
2	Compute an SVD of the small matrix: $\mathbf{B} = \tilde{\mathbf{U}} \mathbf{\Sigma} \mathbf{V}^*$ .
3	Form the orthonormal matrix $\mathbf{U} = \mathbf{Q} \tilde{\mathbf{U}}$ .

Although using *Direct SVD* 2.0.2 algorithm for Stage 2 does not incur additional errors, the computation of  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$  is in general too expensive for dense matrices.

More concretely, the product costs  $\mathcal{O}(mn\ell)$ , even more expensive that the cost of Stage 1 when using Algorithm 1.0.4.

2.0.3. *Faster Procedures.* In order to match the complexity  $\mathcal{O}(mn \log(\ell))$  from Stage 1 we must avoid the product  $\mathbf{Q}^* \mathbf{A}$ .

In [5], the authors propose algorithms based on row extraction of  $\mathbf{Q}$  via its *Interpolative Decomposition*  $\mathbf{Q} = \mathbf{X} \mathbf{Q}_{(J,:)} [2]$ . The proposed algorithm takes  $\mathbf{Q}$  as input and constructs a rank- $k$  matrix factorization

$$(2.1) \quad \mathbf{A} \approx \mathbf{X} \mathbf{B}$$

where  $\mathbf{B}$  is a  $k \times n$  matrix consisting of  $k$  rows extracted from  $\mathbf{A}$ .

The key here is that 2.1 can be produced without any matrix-matrix multiplication resulting in a total of  $\mathcal{O}(k^2(m+n))$  operations. The drawback is that the initial error is larger than the one incurred by  $\mathbf{Q}^* \mathbf{Q} \mathbf{A}$ .

## 3. FULL ALGORITHMS

## THEORY

**3.1. Analysis of Stage 1.** This section focuses on assessing the quality of the basis given by Proto-Algorithm 0.2. More precisely, we want to prove rigorous bounds on the approximation error

$$\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|$$

where  $\|\cdot\|$  denotes either the operator norm or Frobenius norm.

We will split the argument into two parts <sup>3</sup>:

- (1) Provide a generic error bound that depends on the interaction between the test matrix  $\mathbf{\Omega}$  and the right and left singular values of  $\mathbf{A}$ . <sup>4</sup>
- (2) Estimate the error using the distribution of the random matrix. We provide both expectation and probability tail bounds for the error.

**3.1.1. (1) Error bounds via Linear Algebra.** As we aim to compute a rank- $k$  approximation of  $\mathbf{A}$ , we appropriately partition the exact SVD as

$$(3.1) \quad \mathbf{A} = \mathbf{U} \begin{bmatrix} k & n-k \\ \mathbf{\Sigma}_1 & \mathbf{\Sigma}_2 \end{bmatrix} \begin{bmatrix} n \\ \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix}$$

Now, let  $\mathbf{\Omega}_i = \mathbf{V}_i\mathbf{\Omega}$  for  $i = 1, 2$ . Express  $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$  as

$$\mathbf{Y} = \mathbf{A}\mathbf{\Omega} = \mathbf{U} \begin{bmatrix} \ell \\ \mathbf{\Sigma}_1\mathbf{\Omega}_1 \\ \mathbf{\Sigma}_2\mathbf{\Omega}_2 \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix}$$

where  $\mathbf{\Sigma}_1\mathbf{\Omega}_1$  controls most of the action of  $\mathbf{Y}$ , and  $\mathbf{\Sigma}_2\mathbf{\Omega}_2$  is a small perturbation.

The Proto-Algorithm 0.2 computes an orthogonal basis  $\mathbf{Q}$  of  $\text{Im}(\mathbf{Y})$ . In other words, we can express the orthogonal projection to  $\text{Im}(\mathbf{Y})$  as  $\mathbf{P}_\mathbf{Y} = \mathbf{P}_{\text{Im}(\mathbf{Y})} = \mathbf{Q}\mathbf{Q}^*$  <sup>5</sup>. The following Theorem 3.1 bounds the squared error provides a deterministic error bound to the squared error.

**Theorem 3.1** (Deterministic error bound). *We have that*

$$(3.2) \quad \|\mathbf{(\mathbf{I} - \mathbf{P}_\mathbf{Y})\mathbf{A}}\|^2 \leq \|\mathbf{\Sigma}_2\|^2 + \|\mathbf{\Sigma}_2\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|^2,$$

where  $\|\cdot\|$  denotes either the spectral norm or the Frobenius norm.

*Remark 3.2.* Note that  $\mathbf{\Sigma}_1$  does not appear in the error bound. EXPLAIN

*Remark 3.3.* The first term is a deterministic clean error term; we want to compute a rank- $k$  approximation so the error can not be smaller than this term. The second term is a random term that depends on the interaction of the right singular values of  $\mathbf{A}$  amplified by  $\mathbf{\Sigma}_2$ .

We would also like to be able to analyze the power scheme described in REF, i.e,  $\mathbf{B} = (\mathbf{A}\mathbf{A}^*)\mathbf{A} = \mathbf{U}\mathbf{\Sigma}^{2q+1}\mathbf{V}^*$ . The rationale behind the power scheme was that the random approximation of the  $k$ -dimensional gross action of  $\mathbf{A}$  can be improved if we amplify  $\mathbf{\Sigma}_1 - \mathbf{\Sigma}_2$  by power iteration. This can be easily verified by this simple Theorem 3.4.

**Theorem 3.4** (Power scheme). *Let  $\mathbf{A}$  be an  $m \times n$  matrix, and let  $\mathbf{\Omega}$  be an  $n \times \ell$  matrix. Fix a nonnegative integer  $q$ , form  $\mathbf{B} = (\mathbf{A}^*\mathbf{A})^q\mathbf{A}$ , and compute the sample matrix  $\mathbf{Z} = \mathbf{B}\mathbf{\Omega}$ . Then*

$$\|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\mathbf{A}\| \leq \|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\mathbf{B}\|^{1/(2q+1)}.$$

*Remark 3.5.* Let's consider the operator norm, i.e,  $\|\mathbf{\Sigma}_1\| = \sigma_{k+1}$ . Then

$$\|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\mathbf{A}\| \leq \|(\mathbf{I} - \mathbf{P}_\mathbf{Z})\mathbf{B}\|^{1/(2q+1)} \leq \left(1 + \|\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|^2\right)^{1/(4q+2)} \sigma_{k+1}$$

so the power scheme shrinks the suboptimality exponentially fast.

<sup>3</sup>The authors argue that this bipartite proof is common in the literature of randomized linear algebra

<sup>4</sup>Note that we do not deal with randomness yet.

<sup>5</sup>We simplify the notation of the orthogonal projectoin to  $\mathbf{P}_\mathbf{Y}$



Finally, we can ask what are the consequences of truncating the SVD of  $\mathbf{P_Z A}$ , i.e., compute its best rank- $k$  approximation.

**Theorem 3.6** (Analysis of Truncated SVD). *Let  $\mathbf{A}$  be an  $m \times n$  matrix with singular values  $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots$ , and let  $\mathbf{Z}$  be an  $m \times \ell$  matrix, where  $\ell \geq k$ . Suppose that  $\hat{\mathbf{A}}_{(k)}$  is a best rank- $k$  approximation of  $\mathbf{P_Z A}$  with respect to the spectral norm. Then*

$$\|\mathbf{A} - \hat{\mathbf{A}}_{(k)}\| \leq \sigma_{k+1} + \|(\mathbf{I} - \mathbf{P_Z})\mathbf{A}\|.$$

*Remark 3.7.* The result of Theorem 3.6 is quite pessimistic, and in practice we observe that truncating the SVD is not that damaging in the randomized setting.

3.1.2. (2) *Bounds on the gaussian setting.* First we start by providing a bunch of results on gaussian matrices that will be key to prove the bounds on expectation and probability tails.

**Proposition 3.8** (Expected norm of a scaled Gaussian matrix). *Fix matrices  $\mathbf{S}, \mathbf{T}$ , and draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$(3.3) \quad \left(\mathbb{E} \|\mathbf{SGT}\|_{\text{F}}^2\right)^{1/2} = \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|_{\text{F}} \quad \text{and} \quad \mathbb{E} \|\mathbf{SGT}\| \leq \|\mathbf{S}\| \|\mathbf{T}\|_{\text{F}} + \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|.$$

**Proposition 3.9** (Expected norm of a pseudo-inverted Gaussian matrix). *Draw a  $k \times (k + p)$  standard Gaussian matrix  $\mathbf{G}$  with  $k \geq 2$  and  $p \geq 2$ . Then*

$$(3.4) \quad \left(\mathbb{E} \|\mathbf{G}^\dagger\|_{\text{F}}^2\right)^{1/2} = \sqrt{\frac{k}{p-1}} \quad \text{and} \quad \mathbb{E} \|\mathbf{G}^\dagger\| \leq \frac{e\sqrt{k+p}}{p}.$$

**Proposition 3.10** (Concentration for functions of a Gaussian matrix). *Suppose that  $h$  is a Lipschitz function on matrices:*

$$|h(\mathbf{X}) - h(\mathbf{Y})| \leq L \|\mathbf{X} - \mathbf{Y}\|_{\text{F}} \quad \text{for all } \mathbf{X}, \mathbf{Y}.$$

*Draw a standard Gaussian matrix  $\mathbf{G}$ . Then*

$$\mathbb{P}\{h(\mathbf{G}) \geq \mathbb{E} h(\mathbf{G}) + Lt\} \leq e^{-t^2/2}.$$

Now, we are ready to state and prove the main theorems in expectations, and afterwards we will confirm that the error does not oscillate too much around the mean by proving the corresponding bounds on the tails of the distribution.

**Theorem 3.11** (Average Frobenius error). *The expected approximation error can be bounded as follows*

(1)

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P_Y})\mathbf{A}\|_{\text{F}} \leq \sigma_{k+1} \sqrt{\left(1 + \frac{k}{p-1}\right) r(\mathbf{\Sigma_2})}$$

(2)

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P_Y})\mathbf{A}\| \leq \sigma_{k+1} \left(1 + \sqrt{\frac{k}{p-1}} + \frac{e\sqrt{k+p}}{p} \sqrt{r(\mathbf{\Sigma_2})}\right)$$

One interesting quantity is worth examining is  $\|(\mathbf{I} - \mathbf{P_Y})\mathbf{A}\| / \sigma_{k+1}$  to check the factor that tells how far the approximation is from the optimal rank- $k$  approximation 0.5. We observe that the suboptimality term increases essentially as  $\sim \sqrt{k/p}$  and has a term corresponding to the numerical rank 0.2 of the singular values corresponding to the perturbation.<sup>6</sup>

*Proof.* Hölder's inequality and Theorem 3.1 give

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P_Y})\mathbf{A}\|_{\text{F}} \leq \left(\mathbb{E} \|(\mathbf{I} - \mathbf{P_Y})\mathbf{A}\|_{\text{F}}^2\right)^{1/2} \leq \left(\|\mathbf{\Sigma_2}\|_{\text{F}}^2 + \mathbb{E} \|\mathbf{\Sigma_2 \Omega_2 \Omega_1^\dagger}\|_{\text{F}}^2\right)^{1/2}.$$

Then, we condition on  $\mathbf{\Omega_1}$  and use Proposition 3.8 and first part of Proposition 3.9

$$\mathbb{E} \|\mathbf{\Sigma_2 \Omega_2 \Omega_1^\dagger}\|_{\text{F}}^2 = \mathbb{E} \left(\mathbb{E} \left[\|\mathbf{\Sigma_2 \Omega_2 \Omega_1^\dagger}\|_{\text{F}}^2 \mid \mathbf{\Omega_1}\right]\right) = \mathbb{E} \left(\|\mathbf{\Sigma_2}\|_{\text{F}}^2 \|\mathbf{\Omega_1^\dagger}\|_{\text{F}}^2\right) = \|\mathbf{\Sigma_2}\|_{\text{F}}^2 \cdot \mathbb{E} \|\mathbf{\Omega_1^\dagger}\|_{\text{F}}^2 = \frac{k}{p-1} \cdot \|\mathbf{\Sigma_2}\|_{\text{F}}^2,$$

<sup>6</sup>The original statement of Theorem 3.11 from [5] does not explicitly write  $r(\mathbf{\Sigma_2})$ . However, given that we introduced the concept of numerical rank 0.2 and its interpretation, I found interesting to highlight its appearance in the theorem.

Putting everything together

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|_F \leq \left(1 + \frac{k}{p-1}\right)^{1/2} \|\Sigma_2\|_F.$$

and the first part is proved.

The bound on the operator norm is very similar, Theorem 3.1 implies that

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \mathbb{E} \left( \|\Sigma_2\|^2 + \|\Sigma_2 \Omega_2 \Omega_1^\dagger\|^2 \right)^{1/2} \leq \|\Sigma_2\| + \mathbb{E} \|\Sigma_2 \Omega_2 \Omega_1^\dagger\|.$$

Conditioning again on  $\Omega_1$ , we can bound the expectation w.r.t  $\Omega_2$

$$\mathbb{E} \|\Sigma_2 \Omega_2 \Omega_1^\dagger\| \leq \mathbb{E} \left( \|\Sigma_2\| \|\Omega_1^\dagger\|_F + \|\Sigma_2\|_F \|\Omega_1^\dagger\| \right) \leq \|\Sigma_2\| \left( \mathbb{E} \|\Omega_1^\dagger\|_F^2 \right)^{1/2} + \|\Sigma_2\|_F \cdot \mathbb{E} \|\Omega_1^\dagger\|.$$

Finally applying Proposition 3.9, we get to the final result

$$\mathbb{E} \|\Sigma_2 \Omega_2 \Omega_1^\dagger\| \leq \sqrt{\frac{k}{p-1}} \|\Sigma_2\| + \frac{e\sqrt{k+p}}{p} \|\Sigma_2\|_F.$$

□

Finally, we will state the bounds on the tails that prove that the previously expectation bounds are representative of the random behavior.

**Theorem 3.12** (Deviation bounds for the Frobenius error). *Frame the hypotheses of Theorem 3.11. Assume further that  $p \geq 4$ . For all  $u, t \geq 1$ ,*

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\|_F \leq \left(1 + t \cdot \sqrt{12k/p}\right) \left(\sum_{j>k} \sigma_j^2\right)^{1/2} + ut \cdot \frac{e\sqrt{k+p}}{p+1} \cdot \sigma_{k+1},$$

with failure probability at most  $5t^{-p} + 2e^{-u^2/2}$ .

**Theorem 3.13** (Deviation bounds for the spectral error). *Frame the hypotheses of Theorem 3.11, and assume further that  $p \geq 4$ . Then*

$$\|(\mathbf{I} - \mathbf{P}_Y)\mathbf{A}\| \leq \left(1 + 8\sqrt{(k+p) \cdot p \log p}\right) \sigma_{k+1} + 3\sqrt{k+p} \left(\sum_{j>k} \sigma_j^2\right)^{1/2},$$

with failure probability at most  $6p^{-p}$ .

Similar bounds can also be proven for the power scheme [5] that give a high probability guarantee for the bound 3.4.

## EXPERIMENTS

In this section, we briefly present our experiments on some of the algorithms presented at the previous sections. Two sets of experiments are presented.

The first tests are on powers of random gaussian matrices of different sizes. The goal of this set of numerics is to check the performance of the algorithms and make sure they are working as expected. We also analyze the sharpness of the bounds dictated by the theory. In particular, the bound on the expectation of the error given by 3.11 and the error estimation procedure which was motivated by Lemma 1.1.

The second set of experiments is ...

**3.2. Details of the implementation.** The experiments have been performed in **Matlab**. The reproducible source code can be found at the following Github repository together with the  $\text{\LaTeX}$  of the report.

**3.3. Gaussian Matrices.** In this set of experiments, we test and analyze the performance of the Algorithms 1.0.1, 1.0.2, 1.0.3 and 1.0.4.

The experiments are performed on powers of gaussian random matrices of the form:

$$(3.5) \quad \mathbf{A} = \frac{1}{\sqrt{m} + \sqrt{n}} \mathbf{G}, \quad G_{ij} \sim N(0, 1)$$

The normalization in 3.5 is to make sure that the norm of  $\mathbf{A}$  is around 1 with high probability <sup>7</sup>.

3.3.1. *Experiments on Randomized Range Finder 1.0.1.*

3.3.2. *Experiments on Randomized Power Iteration 1.0.3.*

3.3.3. *Experiments on Adaptive Randomized Range Finder 1.0.2.* Show the ratio between the tolerance and the actual error.

3.3.4. *Experiments on Fast Randomized Range Finder 1.0.4.*

**3.4. Real Dataset.**

## CONCLUSION

## REFERENCES

- [1] Åke Björck. “Numerics of gram-schmidt orthogonalization”. In: *Linear Algebra and Its Applications* 197 (1994), pp. 297–316.
- [2] Hongwei Cheng et al. “On the compression of low rank matrices”. In: *SIAM Journal on Scientific Computing* 26.4 (2005), pp. 1389–1404.
- [3] Jack Dongarra and Francis Sullivan. “Guest editors introduction: The top 10 algorithms”. In: *Computing in Science & Engineering* 2.1 (2000), pp. 22–23.
- [4] Ming Gu and Stanley C Eisenstat. “Efficient algorithms for computing a strong rank-revealing QR factorization”. In: *SIAM Journal on Scientific Computing* 17.4 (1996), pp. 848–869.
- [5] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. “Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions”. In: *SIAM review* 53.2 (2011), pp. 217–288.
- [6] Leon Mirsky. “Symmetric gauge functions and unitarily invariant norms”. In: *The quarterly journal of mathematics* 11.1 (1960), pp. 50–59.
- [7] GW Stewart. “The decompositional approach to matrix computation”. In: *Computing in Science & Engineering* 2.1 (2000), pp. 50–59.
- [8] Roman Vershynin. *High Dimensional Probability*. 2016.
- [9] Franco Woolfe et al. “A fast randomized algorithm for the approximation of matrices”. In: *Applied and Computational Harmonic Analysis* 25.3 (2008), pp. 335–366.

DEPARTMENT OF MATHEMATICS, ÉCOLE NORMALE SUPÉRIEURE DE CACHAN  
E-mail address: alexnowakvila@gmail.com

---

<sup>7</sup>This is a direct consequence of Sudakov-Fernique’s inequality that compares the supremum of two gaussian processes when one is dominated by the other. More precisely, you have the bound on the expectation of the norm  $\mathbb{E} \|\mathbf{G}\| \leq \sqrt{m} + \sqrt{n}$  and also an accompanying tail bound  $\mathbb{P} \{\|\mathbf{G}\| \geq \sqrt{m} + \sqrt{n} + t\} \leq 2 \exp(-ct^2)$ . Using Gordon’s inequality (generalization of Sudakov-Fernique’s), you can also prove lower bounds on the smallest singular value  $\mathbb{E} \|\mathbf{G}^\dagger\| \geq \sqrt{m} - \sqrt{n}$  and  $\mathbb{P} \{\|\mathbf{G}\| \leq \sqrt{m} - \sqrt{n} - t\} \leq 2 \exp(-ct^2)$ . These results on concentration of measure have been studied at the Theory Reading Group following the book on High Dimensional Probability from R. Vershynin which I highly recommend [8].

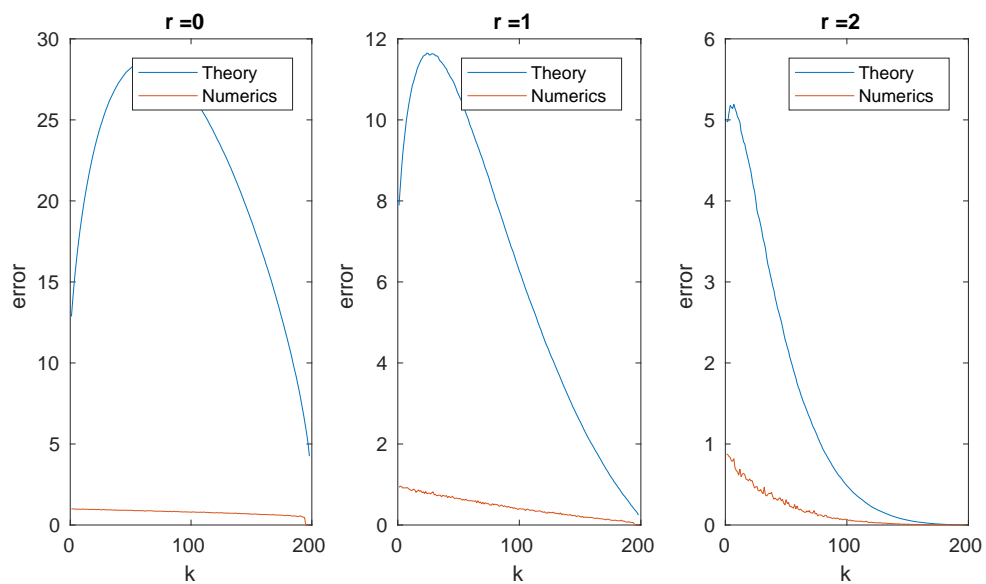


FIGURE 1. bla

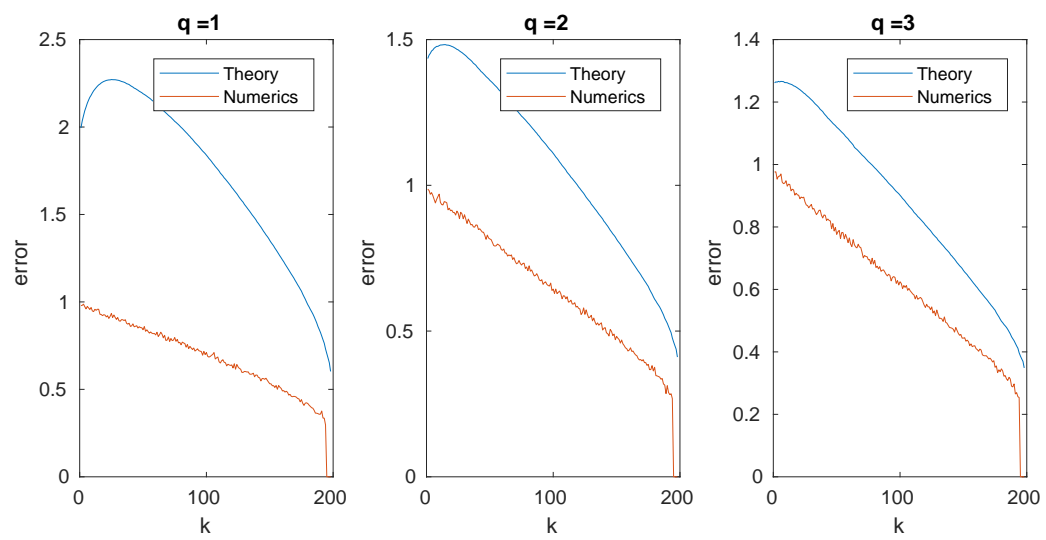


FIGURE 2. bla