

# PROBABILISTIC ALGORITHMS FOR CONSTRUCTING APPROXIMATE 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 randomized methods 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.

## 1. 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 [8] 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.

Data matrices 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  $\mathbf{C} \in \mathbb{R}^{k \times n}$ .

$$(1.1) \quad \begin{matrix} \mathbf{A} & \approx & \mathbf{B} & \mathbf{C}, \\ m \times n & & m \times k & k \times n. \end{matrix}$$

The matrix  $\mathbf{B} \times \mathbf{C}$  in 1.1 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

---

As a project of the course: *Sparsity and Compressed Sensing* at ENS Cachan given by Pr. Gabriel Peyré.

$$(1.2) \quad r(\mathbf{A}) := \frac{\|\mathbf{A}\|_F^2}{\|\mathbf{A}\|^2} = \sum_{j=1}^{\min(m,n)} \left( \frac{\sigma_j}{\sigma_1} \right)^2$$

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 [10] 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

$$(1.3) \quad \mathbf{A} \approx \mathbf{Q}\mathbf{Q}^* \mathbf{A}$$

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 these methods will belong to this first stage.

- **Stage 2:** Given a matrix  $\mathbf{Q}$  that satisfies 1.3, 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.

**1.1. 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 approximation 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

$$(1.4) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^* \mathbf{A}\| \leq \epsilon$$

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 approximation problem*, which seeks to find the best rank- $k$  approximation of the matrix.

$$(1.5) \quad \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\|.$$

---

<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 [10], 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 Singular Value Decomposition (SVD) <sup>2</sup> is key to analyze this problem. Recall that the SVD of a matrix  $\mathbf{A}$  is the following decomposition

$$(1.6) \quad \mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^* = \sum_{j=1}^{\text{rank}(\mathbf{A})} \sigma_j \mathbf{u}_j \mathbf{v}_j^*$$

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.

The SVD provides an optimal answer to the *fixed-rank approximation problem* [7] through the following important observation

$$(1.7) \quad \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\| = \sigma_{k+1}.$$

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 (k-SVD) of the matrix  $\mathbf{A}$ . More precisely, we have that  $\mathbf{B} = \mathbf{U}\mathbf{\Sigma}_{[k]}^{1/2}$  and  $\mathbf{C} = \mathbf{\Sigma}_{[k]}^{1/2}\mathbf{V}^*$  are the best solutions 1.1 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

$$(1.8) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\| \approx \min_{\text{rank}(\mathbf{X}) \leq k} \|\mathbf{A} - \mathbf{X}\|$$

where  $p$  is called the *oversampling parameter*.

**1.2. 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  $\mathbf{w}^{(i)}$  from the input space. Let us repeat this process  $k$  times:

$$(1.9) \quad \mathbf{y}^{(i)} = \mathbf{A}\mathbf{w}^{(i)}, \quad i = 1, \dots, k$$

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  $k$  <sup>3</sup>? 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 1.10, 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:

$$(1.10) \quad \mathbf{y}^{(i)} = \mathbf{A}\mathbf{w}^{(i)} = \mathbf{B}\mathbf{w}^{(i)} + \mathbf{E}\mathbf{w}^{(i)}, \quad i = 1, \dots, k + p$$

<sup>2</sup>The Singular Value Decomposition is a central tool for a data scientist. It is the algorithmic tool used to perform *Principal Component Analysis* (PCA) on a dataset  $\mathbf{X} \in \mathbb{R}^{N \times d}$  of size  $N$  and dimensionality  $d$ . PCA finds linear hidden linear structure in the data by looking at the eigenvectors with largest eigenvalues of the empirical covariance matrix of the data  $\hat{\mathbf{\Sigma}}_N = \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{j=1}^N \mathbf{x}_j \mathbf{x}_j^T$ . These eigenvectors and eigenvalues can also be found by computing the SVD of the data matrix  $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$ , because  $\hat{\mathbf{\Sigma}}_N = \mathbf{V}\mathbf{D}^2\mathbf{V}^T$ . One important remark is that PCA does not suffer from the curse of dimensionality. For example, if we assume sub-gaussianity on the random variable  $X$ , then we have  $\mathbb{E}\|\hat{\mathbf{\Sigma}}_N - \mathbf{\Sigma}\| \leq \varepsilon\|\mathbf{\Sigma}\|$  if we take a sample size of  $N \sim \varepsilon^{-2}d$ , i.e, the number of samples scales linearly with the dimensionality of the dataset. If instead of sub-gaussianity, we relax it to  $\|X\|_2 \leq K(\mathbb{E}\|X\|_2^2)^{1/2}$  we only get an extra log-term. This is one of the main reasons why PCA works!

<sup>3</sup>In practice this is always the case.

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; see for example Theorem 3.11. The theory also shows that  $p$  can be quite small. In practice,  $p = 5$  is enough.

Here in the box below, we illustrate the basic procedure to find  $\mathbf{Q}$  that approximates the range of  $\mathbf{A}$  through randomization. We call it Proto-Algorithm.

| 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 $\mathbf{Y}$ . |

**1.3. 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 directly produce a low rank matrix decomposition  $\mathbf{A} \approx \mathbf{BC}$ . However, there are more involved methods to produce a decomposition from  $\mathbf{Q}$  with computational advantages compared to the direct method.

Once we have  $\mathbf{A} \approx \mathbf{BC}$ , standard matrix decompositions (e.g. SVD, QR) can be easily derived as we show in 2.2.1.

A lot of questions need to be addressed in order to turn these methods into a technology or *off-the-shelf* algorithms. The following sections will be devoted to this.

## 2. 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.

We split this section in three main parts: Stage 1 2.1 and Stage 2 2.2 that will study in depth the algorithmic details of both stages, and Full Algorithms 2.3 in which full algorithms will be presented for a diversity of matrices with certain properties.

**2.1. Stage 1.** In the introductory section, we provided some intuition on the randomized procedure and we developed a general Proto-Algorithm 1.2 to find the range-approximating matrix  $\mathbf{Q}$ . However, Proto-Algorithm 1.2 is very general and can be tuned depending on the problem requirements. The number  $T_{\text{basic}}$  of flops required by Proto-Algorithm 1.2 satisfies

$$(2.1) \quad T_{\text{basic}} \sim \ell n T_{\text{rand}} + \ell T_{\text{mult}} + \ell^2 m$$

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. The last term comes from the orthonormalization procedure of  $\mathbf{Y}$ .

We will now describe some specific realizations of Proto-Algorithm 1.2 that will be intended for problems with different requirements.

**2.1.1. Randomized Range Finder.** This is the most naive and simplest implementation of Proto-Algorithm 1.2. Given an oversampling parameter  $p$ , the *Randomized Range Finder* performs Proto-Algorithm 1.2 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.

The complexity analysis of the Algorithm 2.1.1 gives

$$(2.2) \quad T_{\text{Randomized Range Finder}} \sim \mathcal{O}(mn\ell)$$

## 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  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

this is because generating a gaussian random number is  $\mathcal{O}(1)$  and computing a matrix vector multiplication is  $\mathcal{O}(mn)$ <sup>4</sup>. One important observation is that the complexity  $\mathcal{O}(mn\ell)$  can be prohibitive for massive datasets<sup>5</sup>. A variant of this procedure will be studied in 2.1.4 to get around this problem,

**2.1.2. Adaptive Randomized Range Finder.** One important pitfall of the *Randomized Range Finder* 2.1.1 is that it requires to know in advance the target rank  $k$ . However, if we intend to solve the *fixed-precision approximation 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 2.1.** *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  $\{\omega^{(i)} : i = 1, 2, \dots, r\}$  of standard Gaussian vectors. Then*

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

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

Lemma 2.1 says that we can bound the error with high probability using inexpensive computations in an online manner. The Lemma 2.1 applied to our problem reads

$$(2.3) \quad \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\| \leq 10\sqrt{\frac{2}{\pi}} \max_{i=1, \dots, r} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)\mathbf{A}\omega^{(i)}\|$$

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

The high probability bound 2.3 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

$$(2.4) \quad \|(\mathbf{I} - \mathbf{Q}^{(l)}(\mathbf{Q}^{(l)})^*)\mathbf{A}\| \leq \epsilon.$$

We call *Adaptive Randomized Range Finder* 2.1.2 to the algorithm derived from Lemma 2.1 that solves this problem.

One important question regarding Algorithm 2.1.2 is how good the bound given by Lemma 2.1 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 4.2.

**2.1.3. Randomized Power Iteration.** The *Randomized Range Finder* 2.1.1 algorithm assumes that the singular values of the matrix decay fast. This can be intuitively seen from equation 1.10, where the small singular values interfere with the calculation of the range. This intuition is made precise in Theorem 3.11, where the error of the approximation depends on the  $\sigma_{k+1}$  and  $\sum_{j>k} \sigma_j^2$ .

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.

<sup>4</sup>Note that parallel computation can be extremely helpful here to drastically reduce the effective complexity of the product  $\mathbf{A}\mathbf{\Omega}$ .

<sup>5</sup>Obviously, we can't hope to do better than  $\mathcal{O}(mn)$  which is the size of the matrix we want to approximate. The term that we want to address is the linear term in  $\ell$ . Classical algorithms to compute the  $\ell$ -SVD cost  $\mathcal{O}(mn\ell)$ , by first computing a rank-revealing QR factorization [4] and then manipulating the terms to get to the desired factorization.

ADAPTIVE RANDOMIZED RANGE FINDER

```

1  Draw standard Gaussian vectors  $\omega^{(1)}, \dots, \omega^{(r)}$  of length  $n$ .
2  For  $i = 1, 2, \dots, r$ , compute  $\mathbf{y}^{(i)} = \mathbf{A}\omega^{(i)}$ .
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       $\mathbf{Q}^{(j)} = [\mathbf{Q}^{(j-1)} \quad \mathbf{q}^{(j)}]$ .
10     Draw a standard Gaussian vector  $\omega^{(j+r)}$  of length  $n$ .
11      $\mathbf{y}^{(j+r)} = (\mathbf{I} - \mathbf{Q}^{(j)}(\mathbf{Q}^{(j)})^*)\mathbf{A}\omega^{(j+r)}$ .
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  $\mathbf{Q} = \mathbf{Q}^{(j)}$ .
```

$$(2.5) \quad \sigma_j(\mathbf{B}) = \sigma_j(\mathbf{A})^{2q+1}, \quad j = 1, 2, 3, \dots$$

The *Randomized Power Iteration* 2.1.3 algorithm is the same as the *Randomized Range Finder* 2.1.1 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 2.1.3 is essentially the same because it only requires  $2q + 1$  as many matrix-multiplications as Algorithm 2.1.1 but the number  $q$  is in practice 2, 3 or 4. This can be seen from Corollary 3.14, which shows that the power iteration drives the approximation gap to 1 exponentially fast as  $q$  increases.

**2.1.4. Fast Randomized Range Finder.** A simple inspection to equation 2.1 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.

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

$$(2.6) \quad \mathbf{\Omega} = \sqrt{\frac{n}{\ell}} \mathbf{D} \mathbf{F} \mathbf{R},$$

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  $p, q = 1, 2, \dots, n$
- $\mathbf{R}$  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 [11], 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

$$(2.7) \quad T_{\text{struct}} \sim mn \log(\ell) + \ell^2 n$$

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 (2.6).
- 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  $\mathbf{Y} = \mathbf{Q}\mathbf{R}$ .

**2.2. Stage 2.** The output of the Stage 1 2.1 produces an orthonormal matrix  $\mathbf{Q}$  whose range captures the action of the matrix  $\mathbf{A}$ . The goal of Stage 2 2.2 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  $\|\mathbf{A} - \mathbf{B}\mathbf{C}\| \leq \varepsilon$ . 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  $\mathbf{C} = \mathbf{Q}^* \mathbf{A}$ .

**2.2.1. Compute standard QR and SVD from a general factorization.** Now we will specify how we can compute the standard SVD and QR decompositions from a general low rank decomposition  $\|\mathbf{A} - \mathbf{B}\mathbf{C}\| \leq \varepsilon$  maintaining the tolerance  $\varepsilon$  from Stage 1 2.1.

- *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:*  $\|\mathbf{A} - \mathbf{Q}\mathbf{R}\| \leq \varepsilon$ 
  - (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  $\mathbf{Q} = \mathbf{Q}_1 \mathbf{Q}_2$ .

**2.2.2. Direct SVD.** The procedure described in 2.2.1 to compute the approximate SVD decomposition without sacrificing error, defines what we call the *Direct SVD* 2.2.2.

#### 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* 2.2.2 algorithm for Stage 2 2.2 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 than the cost of Stage 1 2.1 when using the accelerated Algorithm 2.1.4.

### 2.2.3. Different Procedures.

- *Match the cost  $\mathcal{O}(mn \log(\ell))$  of accelerated Stage 1 2.1.4.*

In order to match the complexity  $\mathcal{O}(mn \log(\ell))$  from Stage 1 2.1 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]. Now,  $\mathbf{Q}_{(J,:)}$  is a  $k \times k$  matrix. The proposed algorithm takes  $\mathbf{Q}$  as input and constructs a rank- $k$  matrix factorization

$$(2.8) \quad \mathbf{A} \approx \mathbf{X} \mathbf{B}$$

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

The key here is that 2.8 can be produced without any large 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  $\mathbf{Q}^* \mathbf{Q} \mathbf{A}$  <sup>6</sup>.

- *Single-Pass algorithms*

The previously described algorithms require revisit the matrix  $\mathbf{A}$  multiple times. In [5] they propose single-pass algorithms relying on the observation that all the information you need to compute the decomposition is in the matrices  $\mathbf{\Omega}$ ,  $\mathbf{Y}$  and  $\mathbf{Q}$ . They call the Algorithm *Eigenvalue Decomposition in One-pass*. However, this algorithm too adds additional error to the Stage 1 <sup>7</sup>.

**2.3. Full Algorithms.** Let's now propose a full algorithm depending on the input matrix properties.

- *General Matrices That Fit in Core Memory*

In this case, the appropriate method for Stage 1 is the accelerated Algorithm 2.1.4 that uses structured random matrices. For Stage 2, use the *row-extraction technique* described in 2.2.3. The overall cost  $T_{\text{random}}$  reduces to

$$(2.9) \quad \boxed{T_{\text{random}} \sim mn \log(k) + k^2(m+n)}$$

The overall approximation error satisfies <sup>8</sup>

$$(2.10) \quad \|\mathbf{A} - \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*\| \lesssim n \cdot \sigma_{k+1},$$

- *Matrices for which Matrix-Vector Products are Cheap.*

In this set of matrices we include for example sparse and structured matrices for which the matrix-vector multiplication is substantially cheaper than the general case.

The appropriate method for Stage 1 is Algorithm 2.1.1 with  $p$  constant <sup>9</sup>, or more generally Algorithm 2.1.3 with  $q > 0$ . For Stage 2, we apply the *Direct SVD* 2.2.2. The total cost  $T_{\text{sparse}}$  satisfies

$$(2.11) \quad \boxed{T_{\text{sparse}} = (2q+2)(k+p) T_{\text{mult}} + \mathcal{O}(k^2(m+n))}$$

The overall approximation error satisfies (see Corollary 3.14).

$$(2.12) \quad \|\mathbf{A} - \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*\| \lesssim (kn)^{1/2(2q+1)} \cdot \sigma_{k+1}.$$

<sup>6</sup>In [5], they prove that the error is increased at most by  $(\sqrt{1+4k(n-k)}) \varepsilon$ , which can be large when  $\varepsilon$  is not that small or  $kn$  is large.

<sup>7</sup>Authors in [5] argue that this issue can be addressed with extra oversampling.

<sup>8</sup>This can be derived by combining the error analysis of range approximation with SRFT (which can be found in [5]) and the additional error incurred by the *row-extraction technique*.

<sup>9</sup>If we are dealing with the *fixed-precision problem*, then we use Algorithm 2.1.2.



## 3. THEORY

**3.1. Analysis of Stage 1 2.1.** This section focuses on assessing the quality of the basis given by Proto-Algorithm 1.2. More precisely, we want to prove rigorous bounds on the approximation error

$$(3.1) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^*\mathbf{A}\|$$

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

We will split the argument into two parts <sup>10</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>11</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

$$(3.2) \quad \mathbf{A} = \mathbf{U} \begin{bmatrix} k & n-k \\ \mathbf{\Sigma}_1 & \mathbf{\Sigma}_2 \end{bmatrix} \begin{bmatrix} n \\ \mathbf{V}_1^* \\ \mathbf{V}_2^* \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix}$$

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 1.2 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>12</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*

$$(3.3) \quad \|(\mathbf{I} - \mathbf{P}_\mathbf{Y})\mathbf{A}\|^2 \leq \|\mathbf{\Sigma}_2\|^2 + \|\mathbf{\Sigma}_2\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|^2,$$

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.

**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 2.1.3, 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 the following 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>10</sup>The authors argue that this bipartite proof is common in the literature of randomized linear algebra

<sup>11</sup>Note that we do not deal with randomness yet.

<sup>12</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*

$$(3.4) \quad \left( \mathbb{E} \|\mathbf{SGT}\|_{\text{F}}^2 \right)^{1/2} = \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|_{\text{F}} \quad \text{and} \quad \mathbb{E} \|\mathbf{SGT}\| \leq \|\mathbf{S}\| \|\mathbf{T}\|_{\text{F}} + \|\mathbf{S}\|_{\text{F}} \|\mathbf{T}\|.$$

**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*

$$(3.5) \quad \left( \mathbb{E} \|\mathbf{G}^\dagger\|_{\text{F}}^2 \right)^{1/2} = \sqrt{\frac{k}{p-1}} \quad \text{and} \quad \mathbb{E} \|\mathbf{G}^\dagger\| \leq \frac{e\sqrt{k+p}}{p}.$$

**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 proof 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 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 1.5. We observe that the suboptimality term increases essentially as  $\sim \sqrt{k/p}$  and has a term corresponding to the numerical rank 1.2 of the singular values corresponding to the perturbation.<sup>13</sup> We will go through the proof for the sake of illustration.

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

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

<sup>13</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 1.2 and its interpretation, I found interesting to highlight its appearance in the theorem.

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

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

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 if 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.

Finally, let's analyze the error bounds on the the Power Scheme 2.1.3 that are directly derived from Theorem 3.4 using Hölder's inequality and the bound  $\|\Sigma_2^{2q+1}\|_F \leq \left( \sqrt{\min\{m, n\} - k} \right) \sigma_{k+1}^{1/(2q+1)}$ .

**Corollary 3.14** (Average spectral error for the power scheme). *Frame the hypotheses of Theorem 3.11. Define  $\mathbf{B} = (\mathbf{A}\mathbf{A}^*)^q \mathbf{A}$  for a nonnegative integer  $q$ , and construct the sample matrix  $\mathbf{Z} = \mathbf{B}\Omega$ . Then*

$$\mathbb{E} \|(\mathbf{I} - \mathbf{P}_Z) \mathbf{A}\| \leq \left[ 1 + \sqrt{\frac{k}{p-1}} + \frac{e\sqrt{k+p}}{p} \cdot \sqrt{\min\{m, n\} - k} \right]^{1/(2q+1)} \sigma_{k+1}.$$

From Corollary 3.14 we observe how as we increase  $q$ , the power scheme drives the extra factor in the error to one exponentially fast.

## 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 2.1. We have used an *oversampling parameter*  $p = 5$  for all of the experiments.

The second set of experiments analyses the performance of the *Randomized Power Method* 2.1.3 on MNIST and on a matrix appearing in image processing.

**4.1. Details of the implementation.** The experiments have been implemented in `Matlab`. The reproducible source code can be found at the following Github repository together with the `LATEX` of the report. Please read `README` documentation to run the code on the `Matlab` interactive command line.

**4.2. Gaussian Matrices.** In this set of experiments, we test and analyze the performance of the Algorithms 2.1.1, 2.1.2, 2.1.3 and 2.1.4.

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

$$(4.1) \quad \mathbf{A} = \frac{1}{\sqrt{m} + \sqrt{n}} \mathbf{G}, \quad G_{ij} \sim N(0, 1)$$

The normalization in 4.1 is to make sure that the norm of  $\mathbf{A}$  is around 1 with high probability<sup>14</sup>. As seen from Equation 2.5, the singular values decay faster for higher powers of the matrix.

4.2.1. *Experiments on Randomized Range Finder 2.1.1.* See Figure 4.2.1.

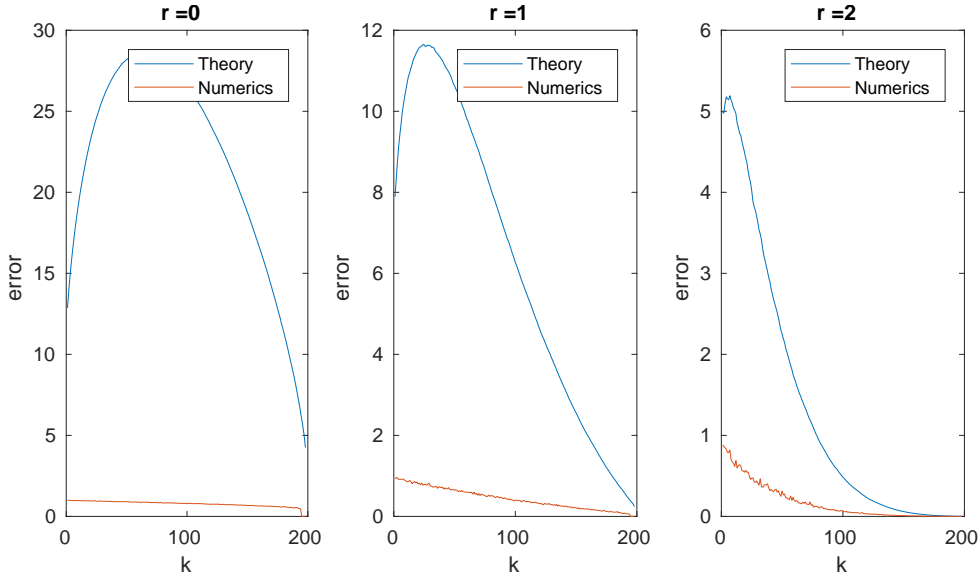


FIGURE 1. Comparison between theoretical mean bound 3.11 and numerical error 3.1 produced by Algorithm 2.1.1 for  $(\mathbf{A}\mathbf{A}^*)^r \mathbf{A}$  with  $r = 0, 1, 2$ .

<sup>14</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 [10].

4.2.2. *Experiments on Randomized Power Iteration 2.1.3.* See Figure 4.2.2.

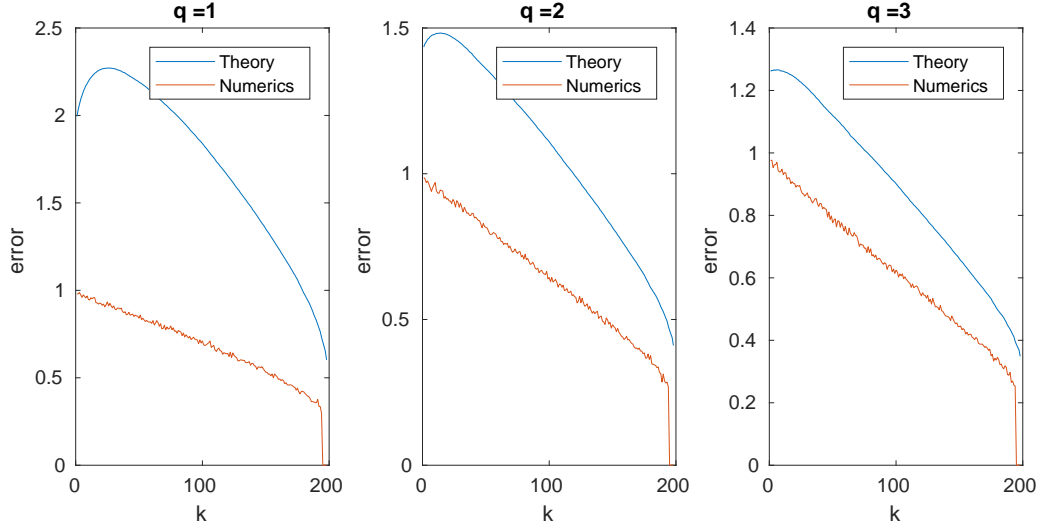


FIGURE 2. Comparison between theoretical mean bound from Corollary 3.14 and numerical error 3.1 produced by Algorithm 2.1.2 for  $\mathbf{A}$  and  $q = 1, 2, 3$

4.2.3. *Experiments on Adaptive Randomized Range Finder 2.1.2.* In this experiment, Algorithm 2.1.2 is applied to  $(\mathbf{A}\mathbf{A}^*)^2\mathbf{A}$  with pre-specified tolerances ranging from 0 to 1. We have fitted a line to the numerical error produced vs the tolerance (the dependence was close to linear). The slope of the resulting line is **0.045**, i.e, for the matrices used in the experiment, the bound 2.3 has a numerical suboptimality of 0.045.

4.2.4. *Experiments on Fast Randomized Range Finder 2.1.4.* Algorithm 2.1.4 has been applied to several powers of  $\mathbf{A}$  (using SRFT test matrices) and the error curve for different target ranks is essentially the same as the one found using Algorithm 2.1.1.

### 4.3. Real Datasets.

4.3.1. *MNIST.* In this experiment, we used 10,000 examples from the MNIST dataset [6] of 28x28 black and white images of handwritten digits.

The goal is to compare Algorithm 2.1.1 with the power version 2.1.3 on the matrix  $\mathbf{M} \in \mathbb{R}^{10000 \times 784}$  where each row is a flattened image. The singular values of  $\mathbf{M}$  decay fast, hence, the power scheme does not significantly improve the randomized scheme for  $q > 1$ . See Figure 4.3.2.

4.3.2. *A Large, Sparse, Noisy Matrix Arising in Image Processing.* One way to tackle standard tasks in image processing is to use the information about the local geometry of the image through the *graph Laplacian*. The dominant eigenvectors of this operator provide coordinates that help to smooth out noisy image patches [9].

First, define  $\tilde{\mathbf{W}}$  as

$$(4.2) \quad \tilde{w}_{ij} = \exp\{-\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2/\sigma^2\}$$

where  $\mathbf{x}^{(i)}$  is a vector gathering the intensities of a neighborhood of a certain size of the pixel  $i$ . Here,  $\sigma$  controls the level of sensitivity. We then define the sparse matrix  $\mathbf{W}$  constructed by zeroing out all the entries in  $\tilde{\mathbf{W}}$  except the 7 largest ones in each row. Finally, the laplacian operator is defined as

$$(4.3) \quad \mathbf{L} = \mathbf{I} - \mathbf{D}^{-1/2}\mathbf{W}\mathbf{D}^{-1/2}$$

where  $\mathbf{D}$  is the diagonal matrix with entries  $d_{ii} = \sum_j w_{ij}$ .

In order to find the low-frequency eigenvectors of  $\mathbf{L}$ , we must find the eigenvectors with dominant eigenvalues of

$$(4.4) \quad \mathbf{A} = \mathbf{D}^{-1/2} \mathbf{W} \mathbf{D}^{-1/2}$$

In our experiment, we have taken a patch of size 57x57 from the classical `Lenna.png` photo. The matrix  $\widehat{\mathbf{W}}$  is constructed with patches of size 5x5 centered at the corresponding pixel with zero-padding at the boundary. The resulting matrix  $\mathbf{A}$  has size 3249x3249 and has low decaying eigenvalues. Hence, as can be seen in Figure 4.3.2, the power scheme from Algorithm 2.1.3 is very advantageous.

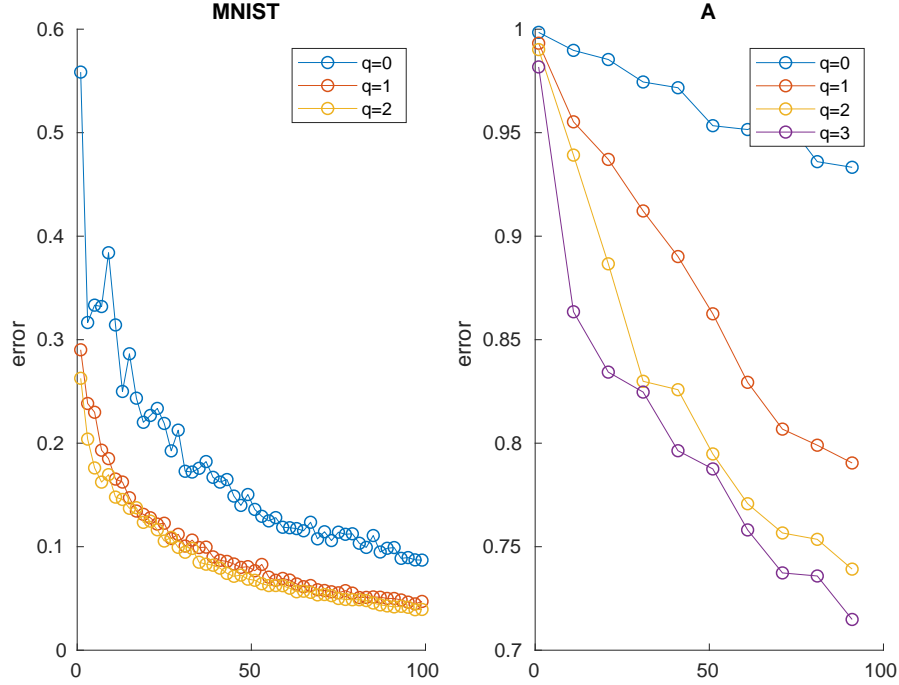


FIGURE 3. Approximation error 3.1 for the power scheme 2.1.3 for MNIST (Left) and  $\mathbf{A}$  4.4 (Right) for  $k$  ranging from 1 to 100.

## 5. CONCLUSION

Randomization schemes for approximate matrix factorization are a powerful and versatile tool that **every data scientist should know about**. These methods are more robust and scalable than classical algorithms, and despite not achieving the same level of accuracy, the error gap can be perfectly controlled with the *oversampling parameter*  $p$  and the number of power iterations  $q$ .

Moreover, these methods rely on theoretical bounds resulting from theorems on concentration of measure. We have seen in Experiment 4.2.1 that the bounds on the expectation can be very loose for certain matrices. One interesting direction of research can be the sharpening of bounds like 3.11 under additional hypothesis on the matrix  $\mathbf{A}$ .

There are several variants in this set of methods that are designed for matrices with different properties, see Section 2.3. One drawback is that sometimes they can be not that easy to verify before applying the algorithm, for e.g, the decaying rate of the spectrum.

From a more personal point of view, I have personally appreciated both the theoretical and experimental section of this project and how they complement each other. I think it is extremely important to have nice experimental results always grounded by theory.

## 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] Yann LeCun. “The MNIST database of handwritten digits”. In: <http://yann.lecun.com/exdb/mnist/> (1998).
- [7] Leon Mirsky. “Symmetric gauge functions and unitarily invariant norms”. In: *The quarterly journal of mathematics* 11.1 (1960), pp. 50–59.
- [8] GW Stewart. “The decompositional approach to matrix computation”. In: *Computing in Science & Engineering* 2.1 (2000), pp. 50–59.
- [9] Arthur D Szlam, Mauro Maggioni, and Ronald R Coifman. “Regularization on graphs with function-adapted diffusion processes”. In: *Journal of Machine Learning Research* 9.Aug (2008), pp. 1711–1739.
- [10] Roman Vershynin. *High Dimensional Probability*. 2016.
- [11] 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.

MASTER “MATHÉMATIQUES, VISION ET APPRENTISSAGE”, ÉCOLE NORMALE SUPÉRIEURE DE CACHAN  
 E-mail address: alexnowakvila@gmail.com