# Chapter 4

# Preliminaries for the investigation of low-rank approximation algorithms

This chapter consolidates probabilistic and linear algebraic tools used in Chapters 5 and 6. We also establish two lemmas of independent interest: the first, Lemma 4.3, is an exponential tail bound on the Frobenius-norm error incurred when approximating the product of two matrices using randomized column and row sampling without replacement; the second, Lemma 4.9, is a deterministic bound on the forward errors of column-based low-rank approximations.

## 4.1 Probabilistic tools

In this section, we review several tools that are used to deal with random matrices and more generally, random processes.

### 4.1.1 Concentration of convex functions of Rademacher variables

Rademacher random variables take the values $\pm 1$ with equal probability. Rademacher vectors are vectors of i.i.d. Rademacher random variables. Rademacher vectors often play a crucial role in the construction of dimension reduction maps, an area where the strong measure concentration properties of Rademacher sums are often exploited. The following result states a large-deviation

property of convex Lipschitz functions of Rademacher vectors: namely, these functions tend to be not much larger than their expectations.

**Lemma 4.1** (A large deviation result for convex Lipschitz functions of Rademacher random variables [Corollary 1.3 ff. in [Led96]] ). *Suppose $f : \mathbb{R}^n \to \mathbb{R}$ is a convex function that satisfies the Lipschitz bound*

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L \left\| \mathbf{x} - \mathbf{y} \right\|_2 \quad \text{for all } \mathbf{x}, \mathbf{y}.$$

*Let $\boldsymbol{\varepsilon} \in \mathbb{R}^n$ be a Rademacher vector. For all $t \geq 0$,*

$$\mathbb{P}\left\{ f(\boldsymbol{\varepsilon}) \geq \mathbb{E}\left[ f(\boldsymbol{\varepsilon}) \right] + Lt \right\} \leq \mathrm{e}^{-t^2/8}.$$

### 4.1.2 Chernoff bounds for sums of random matrices sampled without replacement

Classical Chernoff bounds provide tail bounds for sums of nonnegative random variables. Their matrix analogs provide tail bounds on the eigenvalues and singular values of sums of positive-semidefinite random matrices. Matrix Chernoff bounds are particularly useful for analyzing algorithms involving randomized column-sampling. Most matrix Chernoff bounds available in the literature require the summands to be independent. Indeed, the Chernoff bounds developed in Chapter 2 bound the eigenvalues of a sum of independent random Hermitian matrices. However, occasionally one desires Chernoff bounds that do not require the summands to be independent. The following Chernoff bounds are useful in the case where the summands are drawn without replacement from a set of bounded random matrices.

**Lemma 4.2** (Matrix Chernoff Bounds, Theorem 2.2 in [Tro11b]). *Let $\mathscr{X}$ be a finite set of*

*positive-semidefinite matrices with dimension k, and suppose that*

$$\max_{\mathbf{X} \in \mathcal{X}} \lambda_{\max}(\mathbf{X}) \le B.$$

*Sample $\{\mathbf{X}_1, \ldots, \mathbf{X}_\ell\}$ uniformly at random from $\mathcal{X}$ without replacement. Compute*

$$\mu_{\max} = \ell \cdot \lambda_1(\mathbb{E}\mathbf{X}_1) \quad and \quad \mu_{min} = \ell \cdot \lambda_k(\mathbb{E}\mathbf{X}_1).$$

*Then*

$$\mathbb{P}\left\{\lambda_1\left(\sum_j \mathbf{X}_j\right) \ge (1+v)\mu_{\max}\right\} \le k \cdot \left[\frac{\mathrm{e}^v}{(1+v)^{1+v}}\right]^{\mu_{\max}/B} \quad for\ v \ge 0, and$$

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \le (1-v)\mu_{\min}\right\} \le k \cdot \left[\frac{\mathrm{e}^{-v}}{(1-v)^{1-v}}\right]^{\mu_{\min}/B} \quad for\ v \in [0,1).$$

We also use the following standard simplification of the lower Chernoff bound, which holds under the setup of Lemma 4.2:

$$\mathbb{P}\left\{\lambda_k\left(\sum_j \mathbf{X}_j\right) \le \varepsilon\mu_{\min}\right\} \le k \cdot \mathrm{e}^{-(1-\varepsilon)^2\mu_{\min}/(2B)} \quad for\ \varepsilon \in [0,1]. \tag{4.1.1}$$

### 4.1.3  Frobenius-norm error bounds for matrix multiplication

We now establish a tail bound on the Frobenius-norm error of a simple approximate matrix multiplication scheme based upon randomized column and row sampling. This simple approximate multiplication scheme is a staple in randomized numerical linear algebra, and variants have been analyzed multiple times [DK01, DKM06a, Sar06]. The result derived here differs in that it applies to the sampling without replacement model, and it provides bounds on the error that hold with high probability, rather than simply an estimate of the expected error.

**Lemma 4.3** (Matrix Multiplication). *Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ and $\mathbf{Y} \in \mathbb{R}^{n \times p}$. Fix $\ell \leq n$. Select uniformly at random and without replacement $\ell$ columns from $\mathbf{X}$ and the corresponding rows from $\mathbf{Y}$ and multiply the selected columns and rows with $\sqrt{n/\ell}$. Let $\hat{\mathbf{X}} \in \mathbb{R}^{m \times \ell}$ and $\hat{\mathbf{Y}} \in \mathbb{R}^{\ell \times p}$ contain the scaled columns and rows, respectively. Choose*

$$\sigma^2 \geq \frac{4n}{\ell} \sum\nolimits_{i=1}^{n} \|\mathbf{X}_{(i)}\|_2^2 \|\mathbf{Y}^{(i)}\|_2^2 \quad and \quad B \geq \frac{2n}{\ell} \max_i \|\mathbf{X}_{(i)}\|_2 \|\mathbf{Y}^{(i)}\|_2.$$

*Then if $0 \leq t \leq \sigma^2/B$,*

$$\mathbb{P}\left\{ \left\|\hat{\mathbf{X}}\hat{\mathbf{Y}} - \mathbf{X}\mathbf{Y}\right\|_{\mathrm{F}} \geq t + \sigma \right\} \leq \exp\left(-\frac{t^2}{4\sigma^2}\right).$$

To prove Lemma 4.3, we use the following vector Bernstein inequality for sampling without replacement in Banach spaces; this result follows directly from a similar inequality for sampling with replacement established by Gross in [Gro11]. Again, vector Bernstein inequalities have been derived by multiple authors [LT91, BLM03, Rec11, Tro12, CP11, Gro11]; the value of this specific result is that it applies to the sampling without replacement model.

**Lemma 4.4.** *Let $\mathcal{V}$ be a collection of n vectors in a Hilbert space with norm $\|\cdot\|_2$. Choose $\mathbf{V}_1, \ldots, \mathbf{V}_\ell$ from $\mathcal{V}$ uniformly at random* without *replacement. Choose $\mathbf{V}'_1, \ldots, \mathbf{V}'_\ell$ from $\mathcal{V}$ uniformly at random with* replacement. *Let*

$$\mu = \mathbb{E}\left\|\sum\nolimits_{i=1}^{\ell} (\mathbf{V}'_i - \mathbb{E}\mathbf{V}'_i)\right\|_2$$

*and set*

$$\sigma^2 \geq 4\ell \mathbb{E}\left\|\mathbf{V}'_1\right\|_2^2 \quad and \quad B \geq 2 \max_{\mathbf{V} \in \mathcal{V}} \|\mathbf{V}\|_2.$$

*If $0 \le t \le \sigma^2/B$, then*

$$\mathbb{P}\left\{\left\|\sum_{i=1}^{\ell} \mathbf{V}_i - \ell\mathbb{E}\mathbf{V}_1\right\|_2 \ge \mu + t\right\} \le \exp\left(-\frac{t^2}{4\sigma^2}\right).$$

*Proof.* We proceed by developing a bound on the moment generating function (mgf) of

$$\left\|\sum_{i=1}^{\ell} \mathbf{V}_i - \ell\mathbb{E}\mathbf{V}_1\right\|_2 - \mu.$$

This mgf is controlled by the mgf of a similar sum where the vectors are sampled with replacement. That is, for $\lambda \ge 0$,

$$\mathbb{E}\exp\left(\lambda \cdot \left\|\sum_{i=1}^{\ell} \mathbf{V}_i - \ell\mathbb{E}\mathbf{V}_1\right\|_2 - \lambda\mu\right) \le \mathbb{E}\exp\left(\lambda \cdot \left\|\sum_{i=1}^{\ell} \mathbf{V}_i' - \ell\mathbb{E}\mathbf{V}_1\right\|_2 - \lambda\mu\right). \qquad (4.1.2)$$

This follows from a classical observation due to Hoeffding [Hoe63] that for any convex real-valued function $g$,

$$\mathbb{E}g\left(\sum_{i=1}^{\ell} \mathbf{V}_i\right) \le \mathbb{E}g\left(\sum_{i=1}^{\ell} \mathbf{V}_i'\right).$$

The paper [GN10] provides an alternate exposition of this fact. Specifically, take $g(\mathbf{V}) = \exp\left(\lambda\left\|\mathbf{V} - \ell\mathbb{E}\mathbf{V}_1\right\|_2 - \lambda\mu\right)$ to obtain the inequality of mgfs asserted in (4.1.2).

In the proof of Theorem 12 in [Gro11], Gross establishes that any random variable $Z$ whose mgf is less than the righthand side of (4.1.2) satisfies a tail inequality of the form

$$\mathbb{P}\left\{Z \ge \mu + t\right\} \le \exp\left(-\frac{t^2}{4s^2}\right) \qquad (4.1.3)$$

when $t \le s^2/M$, where

$$s^2 \ge \sum_{i=1}^{\ell} \mathbb{E}\left\|\mathbf{V}_i' - \mathbb{E}\mathbf{V}_1'\right\|_2^2$$

and $\left\|\mathbf{V}_i' - \mathbb{E}\mathbf{V}_1'\right\|_2 \leq M$ almost surely for all $i = 1,\ldots,\ell$. To apply this result, note that for all $i = 1,\ldots,\ell$,

$$\left\|\mathbf{V}_i' - \mathbb{E}\mathbf{V}_1'\right\|_2 \leq 2 \max_{\mathbf{V} \in \mathscr{V}} \|\mathbf{V}\|_2 = B.$$

Take $\mathbf{V}_1''$ to be an i.i.d. copy of $\mathbf{V}_1'$ and observe that, by Jensen's inequality,

$$
\begin{aligned}
\sum_{i=1}^{\ell} \mathbb{E}\left\|\mathbf{V}_i' - \mathbb{E}\mathbf{V}_1'\right\|_2^2 &= \ell\mathbb{E}\left\|\mathbf{V}_1' - \mathbb{E}\mathbf{V}_1'\right\|_2^2 \\
&\leq \ell\mathbb{E}\left\|\mathbf{V}_1' - \mathbf{V}_1''\right\|_2^2 \leq \ell\mathbb{E}\left(\left\|\mathbf{V}_1'\right\|_2 + \left\|\mathbf{V}_1''\right\|_2\right)^2 \\
&\leq 2\ell\mathbb{E}\left\|\mathbf{V}_1'\right\|_2^2 + \left\|\mathbf{V}_1''\right\|_2^2 \\
&= 4\ell\mathbb{E}\left\|\mathbf{V}_1'\right\|_2^2 \leq \sigma^2.
\end{aligned}
$$

The bound given in the statement of Lemma 4.4 when we take $s^2 = \sigma^2$ and $M = B$ in (4.1.3). $\quad\square$

With this Bernstein bound in hand, we proceed to the proof of Lemma 4.3. Let $\mathrm{vec}: \mathbb{R}^{m\times n} \to \mathbb{R}^{mn}$ denote the operation of vectorization, which stacks the columns of a matrix $\mathbf{A} \in \mathbb{R}^{m\times n}$ to form the vector $\mathrm{vec}(\mathbf{A})$.

*Proof of Lemma 4.3.* Let $\mathscr{V}$ be the collection of vectorized rank-one products of columns of $\sqrt{n/\ell} \cdot \mathbf{X}$ and rows of $\sqrt{n/\ell} \cdot \mathbf{Y}$. That is, take

$$\mathscr{V} = \left\{\frac{n}{\ell}\mathrm{vec}(\mathbf{X}_{(i)}\mathbf{Y}^{(i)})\right\}_{i=1}^{n}.$$

Sample $\mathbf{V}_1,\ldots,\mathbf{V}_\ell$ uniformly at random from $\mathscr{V}$ without replacement, and observe that $\mathbb{E}\mathbf{V}_i = \ell^{-1}\mathrm{vec}(\mathbf{XY})$. With this notation,

$$\left\|\hat{\mathbf{X}}\hat{\mathbf{Y}} - \mathbf{XY}\right\|_{\mathrm{F}} \sim \left\|\sum_{i=1}^{\ell}(\mathbf{V}_i - \mathbb{E}\mathbf{V}_i)\right\|_2,$$

where $\sim$ refers to identical distributions. Therefore any probabilistic bound developed for the right-hand side quantity holds for the left-hand side quantity. The conclusion of the lemma follows when we apply Lemma 4.4 to bound the right-hand side quantity.

We calculate the variance-like term in Lemma 4.4,

$$4\ell\mathbb{E}\|\mathbf{V}_1\|_2^2 = 4\ell\frac{1}{n}\sum_{i=1}^n \frac{n^2}{\ell^2}\|\mathbf{X}_{(i)}\|_2^2\|\mathbf{Y}^{(i)}\|_2^2 = 4\frac{n}{\ell}\sum_{i=1}^n \|\mathbf{X}_{(i)}\|_2^2\|\mathbf{Y}^{(i)}\|_2^2 \leq \sigma^2.$$

Now we consider the expectation

$$\mu = \mathbb{E}\left\|\sum_{i=1}^\ell (\mathbf{V}_i' - \mathbb{E}\mathbf{V}_i')\right\|_2.$$

In doing so, we will use the notation $\mathbb{E}\left[C \mid A, B, \dots\right]$ to denote the conditional expectation of a random variable $C$ with respect to the random variables $A, B, \dots$. Recall that a Rademacher vector is a random vector whose entries are independent and take the values $\pm 1$ with equal probability. Let $\boldsymbol{\varepsilon}$ be a Rademacher vector of length $\ell$ and sample $\mathbf{V}_1', \dots, \mathbf{V}_\ell'$ and $\mathbf{V}_1'', \dots, \mathbf{V}_\ell''$ uniformly at random from $\mathcal{V}$ with replacement. Now $\mu$ can be bounded as follows:

$$
\begin{aligned}
\mu &= \mathbb{E}\left\|\sum_{i=1}^\ell (\mathbf{V}_i' - \mathbb{E}\mathbf{V}_i')\right\|_2 \\
&\leq \mathbb{E}\left[\left\|\sum_{i=1}^\ell (\mathbf{V}_i' - \mathbf{V}_i'')\right\|_2 \mid \{\mathbf{V}_i'\}, \{\mathbf{V}_i''\}\right] \\
&= \mathbb{E}\left[\left\|\sum_{i=1}^\ell \varepsilon_i(\mathbf{V}_i' - \mathbf{V}_i'')\right\|_2 \mid \{\mathbf{V}_i'\}, \{\mathbf{V}_i''\}, \boldsymbol{\varepsilon}\right] \\
&\leq 2\mathbb{E}\left[\left\|\sum_{i=1}^\ell \varepsilon_i\mathbf{V}_i'\right\|_2 \mid \{\mathbf{V}_i'\}, \boldsymbol{\varepsilon}\right] \\
&\leq 2\sqrt{\mathbb{E}\left[\left\|\sum_{i=1}^\ell \varepsilon_i\mathbf{V}_i'\right\|_2^2 \mid \{\mathbf{V}_i'\}, \boldsymbol{\varepsilon}\right]} \\
&= 2\sqrt{\mathbb{E}\left[\mathbb{E}\left[\sum_{i,j=1}^\ell \varepsilon_i\varepsilon_j\mathbf{V}_i'^T\mathbf{V}_j' \mid \boldsymbol{\varepsilon}\right] \mid \{\mathbf{V}_i'\}\right]} \\
&= 2\sqrt{\mathbb{E}\sum_{i=1}^\ell \|\mathbf{V}_i'\|_2^2}.
\end{aligned}
$$

The first inequality is Jensen's, and the following equality holds because the components of the sequence $\{\mathbf{V}_i' - \mathbf{V}_i''\}$ are symmetric and independent. The next two manipulations are the triangle inequality and Jensen's inequality. This stage of the estimate is concluded by conditioning and using the orthogonality of the Rademacher variables. Next, the triangle inequality and the fact that $\mathbb{E}\|\mathbf{V}_1'\|_2^2 = \mathbb{E}\|\mathbf{V}_1\|_2^2$ allow us to further simplify the estimate of $\mu$ :

$$\mu \le 2\sqrt{\mathbb{E}\sum_{i=1}^{\ell}\|\mathbf{V}_i'\|_2^2} = 2\sqrt{\ell\mathbb{E}\|\mathbf{V}_1\|_2^2} \le \sigma.$$

We also calculate the quantity

$$2\max_{\mathbf{V}\in\mathscr{V}}\|\mathbf{V}\|_2 = \frac{2n}{\ell}\max_{i}\|\mathbf{X}_{(i)}\|_2\|\mathbf{Y}^{(i)}\|_2 \le B.$$

The tail bound given in the statement of the lemma follows from applying Lemma 4.4 with our estimates for $B$, $\sigma^2$, and $\mu$. $\qquad\square$

## 4.2   Linear Algebra notation and results

In subsequent chapters, we use the following partitioned compact SVD to state results for rectangular matrices $\mathbf{A}$ with $\text{rank}(\mathbf{A}) = \rho$ :

$$\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T = \begin{matrix} k & \rho-k \\ \left[\begin{matrix} \mathbf{U}_1 & \mathbf{U}_2 \end{matrix}\right] \end{matrix} \begin{matrix} k & \rho-k \\ \left[\begin{matrix} \boldsymbol{\Sigma}_1 & \\ & \boldsymbol{\Sigma}_2 \end{matrix}\right] \end{matrix} \left[\begin{matrix} \mathbf{V}_1^T \\ \mathbf{V}_2^T \end{matrix}\right]. \tag{4.2.1}$$

Here, $\boldsymbol{\Sigma}_1$ contains the $k$ largest singular values of $\mathbf{A}$ and the columns of $\mathbf{U}_1$ and $\mathbf{V}_1$ respectively span top $k$-dimensional left and right singular spaces of $\mathbf{A}$. The matrix $\mathbf{A}_k = \mathbf{U}_1\boldsymbol{\Sigma}_1\mathbf{V}_1^T$ is the optimal rank-$k$ approximation to $\mathbf{A}$, and $\mathbf{A}_{\rho-k} = \mathbf{A} - \mathbf{A}_k = \mathbf{U}_2\boldsymbol{\Sigma}_2\mathbf{V}_2^T$. The Moore-Penrose

pseudoinverse of $\mathbf{A}$ is denoted by $\mathbf{A}^\dagger$.

When $\mathbf{A}$ is a positive-semidefinite matrix, $\mathbf{U} = \mathbf{V}$ and (4.2.1) becomes the following partitioned eigenvalue decomposition:

$$\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{U}^T = \begin{array}{cc} k & \rho-k \\ \left[\begin{array}{cc} & \\ \mathbf{U}_1 & \mathbf{U}_2 \end{array}\right] \end{array} \begin{array}{cc} k & \rho-k \\ \left[\begin{array}{cc} \boldsymbol{\Sigma}_1 & \\ & \boldsymbol{\Sigma}_2 \end{array}\right] \end{array} \left[\begin{array}{c} \mathbf{U}_1^T \\ \mathbf{U}_2^T \end{array}\right]. \tag{4.2.2}$$

The eigenvalues of an $n \times n$ symmetric matrix $\mathbf{A}$ are ordered $\lambda_1(\mathbf{A}) \geq \cdots \geq \lambda_n(\mathbf{A})$.

The orthoprojector onto the column space of a matrix $\mathbf{A}$ is written $\mathbf{P_A}$ and satisfies

$$\mathbf{P_A} = \mathbf{A}\mathbf{A}^\dagger = \mathbf{A}(\mathbf{A}^T\mathbf{A})^\dagger\mathbf{A}^T.$$

Let $\mathscr{S}$ be a $k$-dimensional subspace of $\mathbb{R}^n$ and $\mathbf{P}_{\mathscr{S}}$ denote the projection onto $\mathscr{S}$. Then the *coherence* of $\mathscr{S}$ is

$$\mu(\mathscr{S}) = \frac{n}{k}\max_i(\mathbf{P}_{\mathscr{S}})_{ii}.$$

The coherence of a matrix $\mathbf{U} \in \mathbb{R}^{n \times k}$ with orthonormal columns is the coherence of the subspace $\mathscr{S}$ which it spans:

$$\mu(\mathbf{U}) := \mu(\mathscr{S}) = \frac{n}{k}\max_i(\mathbf{P}_{\mathscr{S}})_{ii} = \frac{n}{k}\max_i(\mathbf{U}\mathbf{U}^T)_{ii}.$$

The $k$th column of the matrix $\mathbf{A}$ is denoted by $\mathbf{A}_{(k)}$; the $j$th row is denoted by $\mathbf{A}^{(j)}$. The vector $e_i$ is the $i$th element of the standard Euclidean basis (whose dimensionality will be clear from the context).

We often compare SPSD matrices using the semidefinite ordering. In this ordering, $\mathbf{A}$ is greater than or equal to $\mathbf{B}$, written $\mathbf{A} \succeq \mathbf{B}$ or $\mathbf{B} \preceq \mathbf{A}$, when $\mathbf{A} - \mathbf{B}$ is positive semidefinite. Each SPSD matrix $\mathbf{A}$ has a unique square root $\mathbf{A}^{1/2}$ that is also SPSD, has the same eigenspaces as

$\mathbf{A}$, and satisfies $\mathbf{A} = \left(\mathbf{A}^{1/2}\right)^2$. The eigenvalues of an SPSD matrix $\mathbf{A}$ are arranged in weakly decreasing order: $\lambda_{\max}(\mathbf{A}) = \lambda_1(\mathbf{A}) \geq \lambda_2(\mathbf{A}) \geq \cdots \geq \lambda_n(\mathbf{A}) = \lambda_{\min}(\mathbf{A})$. Likewise, the singular values of a rectangular matrix $\mathbf{A}$ with rank $\rho$ are ordered $\sigma_{\max}(\mathbf{A}) = \sigma_1(\mathbf{A}) \geq \sigma_2(\mathbf{A}) \geq \cdots \geq \sigma_\rho(\mathbf{A}) = \sigma_{\min}(\mathbf{A})$. The spectral norm of a matrix $\mathbf{B}$ is written $\|\mathbf{B}\|_2$; its Frobenius norm and trace are written $\|\mathbf{B}\|_F$ and $\mathrm{Tr}(\mathbf{B})$, respectively. The notation $\|\cdot\|_\xi$ indicates that an expression holds for both $\xi = 2$ and $\xi = F$.

### 4.2.1 Column-based low-rank approximation

The remainder of this thesis concerns low-rank matrix approximation algorithms: Chapter 5 provides bounds on the approximation errors of low-rank approximations that are formed using fast orthonormal transformations, and Chapter 6 provides bounds on the approximation errors of a class of low-rank approximations to SPSD matrices.

Both of these low-rank approximation schemes are amenable to interpretation as schemes wherein a matrix is projected onto a subspace spanned by some linear combination of its columns. The problem of providing a general framework for studying the error of these projection schemes is well studied [BMD09, HMT11, BDMI11]. The authors of these works have provided a set of so-called *structural* results: deterministic bounds on the spectral and Frobenius-norm approximation errors incurred by these projection schemes. Structural results allow us to relate the errors of low-rank approximations formed using projection schemes to the optimal errors $\|\mathbf{A} - \mathbf{A}_k\|_\xi$ for $\xi = 2, F$.

Before stating the specific structural results that are used in the sequel, we review the necessary background material on low-rank matrix approximations that are restricted to lie within a particular subspace.

#### 4.2.1.1 Matrix Pythagoras and generalized least-squares regression

Lemma 4.5 is the analog of Pythagoras' theorem in the matrix setting. A proof of this lemma can be found in [BDMI11]. Lemma 4.6 is an immediate corollary that generalizes the Eckart–Young theorem.

**Lemma 4.5.** *If* $\mathbf{X}\mathbf{Y}^T = \mathbf{0}$ *or* $\mathbf{X}^T\mathbf{Y} = \mathbf{0}$, *then*

$$\|\mathbf{X}+\mathbf{Y}\|_F^2 = \|\mathbf{X}\|_F^2 + \|\mathbf{Y}\|_F^2$$

*and*

$$\max\{\|\mathbf{X}\|_2^2, \|\mathbf{Y}\|_2^2\} \le \|\mathbf{X}+\mathbf{Y}\|_2^2 \le \|\mathbf{X}\|_2^2 + \|\mathbf{Y}\|_2^2.$$

**Lemma 4.6.** *Given* $\mathbf{A} \in \mathbb{R}^{m \times n}$ *and* $\mathbf{C} \in \mathbb{R}^{m \times \ell}$, *for all* $\mathbf{X} \in \mathbb{R}^{\ell \times n}$

$$\|\mathbf{A} - \mathbf{P}_\mathbf{C}\mathbf{A}\|_\xi^2 \le \|\mathbf{A} - \mathbf{C}\mathbf{X}\|_\xi^2$$

*for both* $\xi = 2$ *and* $\xi = F$.

*Proof.* Write

$$\mathbf{A} - \mathbf{C}\mathbf{X} = (\mathbf{I} - \mathbf{P}_\mathbf{C})\mathbf{A} + (\mathbf{P}_\mathbf{C}\mathbf{A} - \mathbf{C}\mathbf{X})$$

and observe that

$$((\mathbf{I} - \mathbf{P}_\mathbf{C})\mathbf{A})^T (\mathbf{P}_\mathbf{C}\mathbf{A} - \mathbf{C}\mathbf{X}) = \mathbf{0},$$

so by Lemma 4.5,

$$\|\mathbf{A} - \mathbf{C}\mathbf{X}\|_\xi^2 \ge \|(\mathbf{I} - \mathbf{P}_\mathbf{C})\mathbf{A}\|_\xi^2.$$

$\square$

### 4.2.1.2 Low-rank approximations restricted to subspaces

Given $\mathbf{A} \in \mathbb{R}^{m \times n}$; a target rank $k < n$; another matrix $\mathbf{Y} \in \mathbb{R}^{m \times \ell}$, where $\ell > k$; and a choice of norm $\xi$ ($\xi = 2$ or $\xi = F$), we use the notation $\mathbf{\Pi}_{\mathbf{Y},k}^{\xi}(\mathbf{A})$ to refer to the matrix that lies in the column span of $\mathbf{Y}$, has rank $k$ or less, and minimizes the $\xi$-norm error in approximating $\mathbf{A}$. More concisely, $\mathbf{\Pi}_{\mathbf{Y},k}^{\xi}(\mathbf{A}) = \mathbf{YX}^{\xi}$, where

$$\mathbf{X}^{\xi} = \underset{\mathbf{X} \in \mathbb{R}^{\ell \times n}: \text{rank}(\mathbf{X}) \leq k}{\arg\min} \|\mathbf{A} - \mathbf{YX}\|_{\xi}^{2}.$$

The approximation $\mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A})$ can be computed using the following three-step procedure:

1: Orthonormalize the columns of $\mathbf{Y}$ to construct a matrix $\mathbf{Q} \in \mathbb{R}^{m \times \ell}$.

2: Compute $\mathbf{X}_{\text{opt}} = \arg\min_{\mathbf{X} \in \mathbb{R}^{\ell \times n}, \text{ rank}(\mathbf{X}) \leq k} \left\| \mathbf{Q}^{T}\mathbf{A} - \mathbf{X} \right\|_{F}$.

3: Compute and return $\mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A}) = \mathbf{Q}\mathbf{X}_{\text{opt}} \in \mathbb{R}^{m \times n}$.

There does not seem to be a similarly efficient algorithm for computing $\mathbf{\Pi}_{\mathbf{Y},k}^{2}(\mathbf{A})$.

The following result, which appeared as Lemma 18 in [BDMI11], both verifies the claim that this algorithm computes $\mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A})$ and shows that $\mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A})$ is a constant factor approximation to $\mathbf{\Pi}_{\mathbf{Y},k}^{2}(\mathbf{A})$.

**Lemma 4.7.** *[Lemma 18 in [BDMI11]] Given $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{Y} \in \mathbb{R}^{m \times \ell}$, and an integer $k \leq \ell$, the matrix $\mathbf{Q}\mathbf{X}_{\text{opt}} \in \mathbb{R}^{m \times n}$ described above satisfies $\mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A}) = \mathbf{Q}\mathbf{X}_{\text{opt}}$, can be computed in $O(mn\ell + (m + n)\ell^2)$ time, and satisfies*

$$\left\| \mathbf{A} - \mathbf{\Pi}_{\mathbf{Y},k}^{F}(\mathbf{A}) \right\|_{2}^{2} \leq 2 \left\| \mathbf{A} - \mathbf{\Pi}_{\mathbf{Y},k}^{2}(\mathbf{A}) \right\|_{2}^{2}.$$

## 4.2.2 Structural results for low-rank approximation

The following result, which appears as Lemma 7 in [BMD09], provides an upper bound on the residual error of the low-rank matrix approximation obtained via projections onto subspaces. The paper [HMT11] also supplies an equivalent result.

**Lemma 4.8.** *[Lemma 7 in [BMD09]] Let* $\mathbf{A} \in \mathbb{R}^{m \times n}$ *have rank* $\rho$. *Fix* $k$ *satisfying* $0 \leq k \leq \rho$. *Given a matrix* $\mathbf{S} \in \mathbb{R}^{n \times \ell}$, *with* $\ell \geq k$, *construct* $\mathbf{Y} = \mathbf{AS}$. *If* $\mathbf{V}_1^T \mathbf{S}$ *has full row-rank, then, for* $\xi = 2, \mathrm{F}$,

$$\|\mathbf{A} - \mathbf{P_Y A}\|_\xi^2 \leq \|\mathbf{A} - \mathbf{\Pi}_{\mathbf{Y},k}^\xi(\mathbf{A})\|_\xi^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_\xi^2 + \|\mathbf{\Sigma}_2 \mathbf{V}_2^T \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger\|_\xi^2. \tag{4.2.3}$$

In addition to this bound on the residual error, we use the following novel structural bound on the forward errors of low-rank approximants.

**Lemma 4.9.** *Let* $\mathbf{A} \in \mathbb{R}^{m \times n}$ *have rank* $\rho$. *Fix* $k$ *satisfying* $0 \leq k \leq \rho$. *Given a matrix* $\mathbf{S} \in \mathbb{R}^{n \times \ell}$, *where* $\ell \geq k$, *construct* $\mathbf{Y} = \mathbf{AS}$. *If* $\mathbf{V}_1^T \mathbf{S}$ *has full row-rank, then, for* $\xi = 2, \mathrm{F}$,

$$\|\mathbf{A}_k - \mathbf{P_Y A}\|_\xi^2 \leq \|\mathbf{A} - \mathbf{A}_k\|_\xi^2 + \|\mathbf{\Sigma}_2 \mathbf{V}_2^T \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger\|_\xi^2. \tag{4.2.4}$$

*Proof.* Observe that

$$(\mathbf{A}_k - \mathbf{P_Y A}_k)^T (\mathbf{P_Y A}_{\rho-k}) = \mathbf{0},$$

so Lemma 4.5 implies that

$$\|\mathbf{A}_k - \mathbf{P_Y A}\|_\xi^2 = \|\mathbf{A}_k - \mathbf{P_Y A}_k - \mathbf{P_Y A}_{\rho-k}\|_\xi^2 \leq \|\mathbf{A}_k - \mathbf{P_Y A}_k\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2.$$

Applying Lemma 4.6 with $\mathbf{X} = (\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T$, we see that

$$\|\mathbf{A}_k - \mathbf{P}_{\mathbf{Y}}\mathbf{A}\|_\xi^2 \leq \|\mathbf{A}_k - \mathbf{Y}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2$$

$$= \|\mathbf{A}_k - \mathbf{A}_k \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T + \mathbf{A}_{\rho-k}\mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2$$

$$= \|\mathbf{A}_k - \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T + \mathbf{A}_{\rho-k}\mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2.$$

Since $\mathbf{V}_1^T \mathbf{S}$ has full row rank, $(\mathbf{V}_1^T \mathbf{S})(\mathbf{V}_1^T \mathbf{S})^\dagger = \mathbf{I}_k$. Recall that $\mathbf{A}_k = \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T$ and $\mathbf{A}_{\rho-k} = \mathbf{U}_2 \boldsymbol{\Sigma}_2 \mathbf{V}_2^T$.

Consequently, the above inequality reduces neatly to the desired inequality

$$\|\mathbf{A}_k - \mathbf{P}_{\mathbf{Y}}\mathbf{A}\|_\xi^2 \leq \|\mathbf{A}_k - \mathbf{U}_1 \boldsymbol{\Sigma}_1 \mathbf{V}_1^T + \mathbf{A}_{\rho-k}\mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2$$

$$= \|\mathbf{A}_{\rho-k}\mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger \mathbf{V}_1^T\|_\xi^2 + \|\mathbf{A}_{\rho-k}\|_\xi^2$$

$$= \|\mathbf{A} - \mathbf{A}_k\|_\xi^2 + \|\boldsymbol{\Sigma}_2 \mathbf{V}_2^T \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger\|_\xi^2.$$

$\square$

### 4.2.2.1  A geometric interpretation of the sampling interaction matrix

Let $\boldsymbol{\Omega}_1 = \mathbf{V}_1^T \mathbf{S}$ and $\boldsymbol{\Omega}_2 = \mathbf{V}_2^T \mathbf{S}$ denote the interaction of the sampling matrix $\mathbf{S}$ with the top and bottom right-singular spaces of $\mathbf{A}$. It is evident from Lemmas 4.8 and 4.9 that the quality of the low-rank approximations depend upon the norm of the *sampling interaction matrix*

$$\mathbf{V}_2^T \mathbf{S}(\mathbf{V}_1^T \mathbf{S})^\dagger = \boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger.$$

The smaller the spectral norm of the $\boldsymbol{\Omega}_2 \boldsymbol{\Omega}_1^\dagger$ the more effective $\mathbf{S}$ is as a sampling matrix. To give the sampling interaction matrix a geometric interpretation, we first recall the definition of the

sine between the range spaces of two matrices $\mathbf{M}_1$ and $\mathbf{M}_2$ :

$$\sin^2(\mathbf{M}_1, \mathbf{M}_2) = \|(\mathbf{I} - \mathbf{P}_{\mathbf{M}_1})\mathbf{P}_{\mathbf{M}_2}\|_2.$$

Note that this quantity is *not* symmetric: it measures how well the range of $\mathbf{M}_1$ captures that of $\mathbf{M}_2$ [GV96, Chapter 12].

**Lemma 4.10.** *Fix* $\mathbf{A} \in \mathbb{R}^{m \times n}$, *a target rank* $k$, *and* $\mathbf{S} \in \mathbb{R}^{n \times \ell}$ *where* $\ell > k$. *Assume* $\mathbf{S}$ *has orthonormal columns. Define*

$$\Omega_1 = \mathbf{V}_1^T \mathbf{S} \quad \text{and} \quad \Omega_2 = \mathbf{V}_2^T \mathbf{S}.$$

*Then, if* $\Omega_1$ *has full row-rank,*

$$\|\Omega_2 \Omega_1^\dagger\|_2 = \tan^2(\mathbf{S}, \mathbf{V}_1).$$

*Proof.* Since $\mathbf{V}_1$ and $\mathbf{S}$ have orthonormal columns, we see that

$$
\begin{aligned}
\sin^2(\mathbf{S}, \mathbf{V}_1) &= \left\|(\mathbf{I} - \mathbf{S}\mathbf{S}^T)\mathbf{V}_1\mathbf{V}_1^T\right\|_2^2 \\
&= \left\|\mathbf{V}_1^T(\mathbf{I} - \mathbf{S}\mathbf{S}^T)\mathbf{V}_1\right\|_2 \\
&= \left\|\mathbf{I} - \mathbf{V}_1^T\mathbf{S}\mathbf{S}^T\mathbf{V}_1\right\|_2 \\
&= 1 - \lambda_k(\mathbf{V}_1^T\mathbf{S}\mathbf{S}^T\mathbf{V}_1) \\
&= 1 - \|\Omega_1^\dagger\|^{-2}.
\end{aligned}
$$

The second to last equality holds because $\mathbf{V}_1^T\mathbf{S}$ has $k$ rows and we assumed it has full row-rank. Accordingly,

$$\tan^2(\mathbf{S}, \mathbf{V}_1) = \frac{\sin^2(\mathbf{S}, \mathbf{V}_1)}{1 - \sin^2(\mathbf{S}, \mathbf{V}_1)} = \|\Omega_1^\dagger\|_2^2 - 1.$$

Now observe that

$$\|\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|_2^2 = \left\|(\mathbf{S}^T\mathbf{V}_1)^\dagger\mathbf{S}^T\mathbf{V}_2\mathbf{V}_2^T\mathbf{S}(\mathbf{V}_1^T\mathbf{S})^\dagger\right\|_2$$

$$= \left\|(\mathbf{S}^T\mathbf{V}_1)^\dagger(\mathbf{I} - \mathbf{S}^T\mathbf{V}_1\mathbf{V}_1^T\mathbf{S})(\mathbf{V}_1^T\mathbf{S})^\dagger\right\|_2$$

$$= \left\|(\mathbf{S}^T\mathbf{V}_1)^\dagger\right\|_2^2 - 1$$

$$= \tan^2(\mathbf{S}, \mathbf{V}_1).$$

The second to last equality holds because of the fact that, for any matrix $\mathbf{M}$,

$$\left\|\mathbf{M}^\dagger(\mathbf{I} - \mathbf{M}\mathbf{M}^T)(\mathbf{M}^T)^\dagger\right\|_2 = \left\|\mathbf{M}^\dagger\right\|_2^2 - 1;$$

this identity can be established with a routine SVD argument. $\qquad\square$

Thus, when $\mathbf{S}$ has orthonormal columns and $\mathbf{V}_1^T\mathbf{S}$ has full row-rank, $\|\mathbf{\Omega}_2\mathbf{\Omega}_1^\dagger\|_2$ is the tangent of the largest angle between the range of $\mathbf{S}$ and the top right singular space spanned by $\mathbf{V}_1$. If $\mathbf{V}_1^T\mathbf{S}$ does not have full row-rank, then our derivation above shows that $\sin^2(\mathbf{S}, \mathbf{V}_1) = 1$, meaning that there is a vector in the eigenspace spanned by $\mathbf{V}_1$ which has no component in the space spanned by the sketching matrix $\mathbf{S}$.

We note that $\tan(\mathbf{S}, \mathbf{V}_1)$ also arises in the classical bounds on the convergence of the orthogonal iteration algorithm for approximating the top $k$-dimensional singular spaces of a matrix (see, e.g. [GV96, Theorem 8.2.2]).