

Lecture 13 — 02/28/2024

Instructor: Shashanka Ubaru

Scribe: Shourya Pandey

1 Recap

In the last lecture, we introduced iterative methods, which predate sketching-based methods, for low rank approximation of a matrix. Recall the *Power Method* for computing the top singular vector of a matrix:

Algorithm 1: Power Method

Data: $\mathbf{A} \in \mathbb{R}^{n \times d}$, $q \in \mathbb{N}$

```

1  $z_0 \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ 
2  $z_0 \leftarrow \frac{z_0}{\|z_0\|_2}$ 
3 for  $\ell = 1, 2, \dots, q$  do
4    $z_\ell \leftarrow \mathbf{A}^\top (\mathbf{A} z_{\ell-1})$ 
5    $z_\ell \leftarrow \frac{z_\ell}{\|z_\ell\|_2}$ 
6 return  $z_q$ 

```

The following theorems record the guarantee of the power method in the gapped and gapless cases.

Theorem 1 (Power Method, Gapped). *Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \sigma_{\min(n,d)}$ and top singular vector v_1 , and let $\gamma := \frac{\sigma_1 - \sigma_2}{\sigma_1}$. Then, for any $\epsilon, \delta \in (0, 1)$ with $\delta = \exp(-O(d))$, the Power Method (Algorithm 1) with $q = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\gamma}\right)$ satisfies*

$$\|v_1 - z_q\|_2 \leq \epsilon$$

with probability at least $1 - \delta$. Moreover, the algorithm runs in time $O\left(\text{nnz}(\mathbf{A}) \frac{\log(d/\epsilon) + \log(1/\delta)}{\gamma}\right)$.

Theorem 2 (Power Method, Gapless). *Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be a matrix with singular values $\sigma_1 \geq \sigma_2 \geq \dots \sigma_{\min(n,d)}$ and top singular vector v_1 , and let $\gamma := \frac{\sigma_1 - \sigma_2}{\sigma_1}$. Then, for any $\epsilon, \delta \in (0, 1)$ with $\delta = \exp(-O(d))$, the Power Method (Algorithm 1) with $q = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\epsilon}\right)$ satisfies*

$$\|\mathbf{A} - \mathbf{A}z_qz_q^\top\|_F^2 \leq (1 + \epsilon)\|\mathbf{A} - \mathbf{A}v_1v_1^\top\|_F^2$$

with probability at least $1 - \delta$. Moreover, the algorithm runs in time $O\left(\text{nnz}(\mathbf{A}) \frac{\log(d/\epsilon) + \log(1/\delta)}{\epsilon}\right)$.

Note that either of these guarantees implies

$$\|\mathbf{A}z_q\|_2^2 \geq (1 - \epsilon)^2 \sigma_1^2.$$

In the gapped case, we can closely align the vector z_q with the top singular vector v_1 , while in the gapless case the flexibility of the power method extends to aligning z_q with the eigenspace of right eigenvectors with sufficiently large singular values. Finally, we saw a natural extension of the Power Method, called the *Block Power Method*, for computing the top k singular vectors of \mathbf{A} .

2 Krylov Subspaces

To motivate the definition of Krylov subspaces, consider a linear regression problem of the form

$$\min_{x \in \mathbb{R}^d} F(x) := \frac{1}{2} \|\mathbf{C}x - b\|_2^2 = \min_{x \in \mathbb{R}^d} \left(\frac{1}{2} x^\top \mathbf{A}x - x^\top v + \frac{1}{2} \|b\|_2^2 \right),$$

where $\mathbf{C} \in \mathbb{R}^{n \times d}$, $b \in \mathbb{R}^n$, $\mathbf{A} = \mathbf{C}^\top \mathbf{C}$, and $v = \mathbf{A}^\top b$. This is a convex optimization problem, because

$$\nabla^2 F(x) = \mathbf{C}^\top \mathbf{C} \succeq \mathbf{0}.$$

A natural approach is to use a descent algorithm. the gradient is given by $\nabla F(x) = \mathbf{A}x - v$, so we initialize the gradient descent method with x_0 equal to a multiple of v , then after q descent steps the point x_q is in the span of the vectors $v, \mathbf{A}v, \dots, \mathbf{A}^q v$. We call

$$\mathbf{K}_q(\mathbf{A}, v) := \text{Span}(v, \mathbf{A}v, \dots, \mathbf{A}^q v)$$

the *Krylov subspace* of dimension q generated by (\mathbf{A}, v) . We will also use the notation \mathbf{K}_q for the Krylov subspace in case the matrix \mathbf{A} and the vector v is clear from the context. An equivalent definition is

$$\mathbf{K}_q(\mathbf{A}, v) := \{p(\mathbf{A})v \mid p \text{ is a polynomial of degree at most } q\}.$$

Suppose \mathbf{A} has full rank. Then, the optimal solution to the linear regression problem is $x = \mathbf{A}^{-1}v$. Krylov subspace methods try to avoid the $O(nd^{\omega-1} + d^\omega)$ cost of matrix multiplication ($\mathbf{A} = \mathbf{C}^\top \mathbf{C}$) and matrix inversion by approximating \mathbf{A}^{-1} using polynomials in \mathbf{A} .

Remark 1. The definition of \mathbf{K}_q immediately implies that $\mathbf{K}_{q'} \subseteq \mathbf{K}_q$ for $q' \leq q$. Moreover, $\mathbf{K}_q \subseteq \mathbb{R}^d$, which has dimension d . This implies the existence of an index q_1 such that

$$\mathbf{K}_0 \subsetneq \mathbf{K}_1 \subsetneq \dots \subsetneq \mathbf{K}_{q_1} = \mathbf{K}_{q_1+1}.$$

It can be shown that $\mathbf{K}_{q_1} = \mathbf{K}_q$ for all $q \geq q_1$. Consider the minimal polynomial p of degree $1 \leq r \leq d$ such that $p(\mathbf{A}) = \mathbf{0}$. Then, \mathbf{A}^r is expressible as a linear combination of the matrices $\mathbf{I}, \mathbf{A}, \dots, \mathbf{A}^{r-1}$, which implies that $\mathbf{K}_r = \mathbf{K}_{r-1}$. Therefore, $q_1 \leq r - 1$. A partial converse is also true: there exists a vector z_1 such that q_1 achieves the value $r - 1$.

3 Lanczos Algorithm

Reconsider the problem of finding the top eigenvector of a symmetric matrix \mathbf{A} . The Krylov iteration methods introduced for linear regression apply more generally via a strategy known as *Lanczos algorithm* or *Lanczos iteration*. The Lanczos algorithm takes a symmetric matrix \mathbf{A} and finds a matrix \mathbf{Z}_q which is an orthonormal basis of a certain Krylov subspace $\mathbf{K}(\mathbf{A}, z_1)$, and such that $\mathbf{T}_q := \mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q$ is a tridiagonal matrix. While the eigenvectors and eigenvalues are not apparent from the tridiagonal form, computing \mathbf{T}_q is already a significant step towards it.

Algorithm 2: Lanczos Algorithm

Data: $\mathbf{A} \in \mathbb{R}^{n \times d}, q \in \mathbb{N}$

- 1 $z_0 \leftarrow \mathbf{0}, \beta_1 \leftarrow 0$
 - 2 Choose a starting vector $z_1 \in \mathbb{R}^d$ with unit norm.
 - 3 $z_0 \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$
 - 4 **for** $\ell = 1, 2, \dots, q - 1$ **do**
 - 5 $y_\ell \leftarrow \mathbf{A}z_\ell - \beta_\ell z_{\ell-1}$
 - 6 $\alpha_\ell \leftarrow \langle y_\ell, z_\ell \rangle$
 - 7 $y_\ell \leftarrow y_\ell - \alpha_\ell z_\ell$
 - 8 $\beta_{\ell+1} \leftarrow \|y_\ell\|_2$. If $\beta_{\ell+1} = 0$ then exit the loop.
 - 9 $z_{\ell+1} \leftarrow \frac{y_\ell}{\beta_{\ell+1}}$
 - 10 $\mathbf{Z}_q \leftarrow [z_1 \ z_2 \ \dots \ z_q]$
 - 11 **return** \mathbf{Z}_q
-

The matrix $\mathbf{T}_q = \mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q$ is called the *Rayleigh Ritz-projection* and is given by

$$\mathbf{T}_q = \begin{bmatrix} \alpha_1 & \beta_2 & & & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & & & \\ & \beta_3 & \alpha_3 & \beta_4 & & & & \\ & & & \cdot & \cdot & \cdot & & \\ & & & & \cdot & \cdot & \cdot & \\ & & & & & \beta_q & \alpha_q & \end{bmatrix}$$

If u is a top eigenvector estimate of \mathbf{T}_q , then $\mathbf{Z}_q u$ is the estimate of the eigenvector of \mathbf{A} .

Theorem 3 (Lanczos Algorithm, Gapped). *Let $\gamma := \frac{\lambda_1 - \lambda_2}{\lambda_1}$ be the gap between the largest eigenvalue, λ_1 , and the second largest eigenvalue, λ_2 , of $\mathbf{A} \in \mathbb{S}_{\neq \mathbf{0}}^{d \times d}$, and let v_1 be the top eigenvector of \mathbf{A} ¹. Let $\epsilon, \delta \in (0, 1)$ with $\delta = \exp(-O(d))$. If the Lanczos's algorithm (Algorithm 2) is initialized with a normalized random Gaussian vector with $q = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\gamma}}\right)$, and u is the top eigenvector of $\mathbf{T}_q = \mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q$, then the vector $w = \mathbf{Z}_q u$ satisfies*

$$\|\mathbf{A} - \mathbf{A} w w^\top\|_F^2 \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A} v_1 v_1^\top\|_F^2 \quad (1)$$

with probability at least $1 - \delta$. Moreover, the algorithm takes time $O\left(\text{nnz}(\mathbf{A}) \frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\gamma}}\right)$ ².

Proof. First, assuming \mathbf{Z}_q has full rank, we claim that the amongst all vectors that span the Krylov subspace $\mathbf{K}_q(\mathbf{A}, z_1)$ (which is also the span of the columns of \mathbf{Z}_q), the vector $w = \mathbf{Z}_q u$ minimizes $\|\mathbf{A} - \mathbf{A} w w^\top\|_F^2$. Any vector in the span of \mathbf{Z}_q is of the form $y = \mathbf{Z}_q x$ for some $x \in \mathbb{R}^q$. Now,

$$\begin{aligned} \|\mathbf{A} - \mathbf{A} y y^\top\|_F^2 &= \|\mathbf{A} - \mathbf{A} \mathbf{Z}_q x x^\top \mathbf{Z}_q^\top\|_F^2 \\ &= \text{Tr} \left(\mathbf{A}^\top \mathbf{A} - \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x x^\top \mathbf{Z}_q^\top - \mathbf{Z}_q x x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} + \mathbf{Z}_q x x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x x^\top \mathbf{Z}_q^\top \right) \\ &= \text{Tr} \left(\mathbf{A}^\top \mathbf{A} - 2x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x + \left(x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x\right) \left(x^\top \mathbf{Z}_q^\top \mathbf{Z}_q x\right) \right) \\ &= \text{Tr} \left(\mathbf{A}^\top \mathbf{A} \right) - 2x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x + \left(x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x\right) \|x\|_2^2. \end{aligned}$$

¹Variants of the Lanczos algorithm work for non-symmetric matrices too, such as Arnoldi's iterations.

²The time taken to compute the top eigenvector u of \mathbf{T}_q is $O(q^3)$ and can be made as small as $O(q \log q)$ via the Fast Multipole Method [1] for tridiagonal matrices.

The second equality used $\|\mathbf{X}\|_F^2 = \text{Tr}(\mathbf{X}^\top \mathbf{X})$, the third equality used the cyclic property of trace, and the fourth equality used $\mathbf{Z}_q^\top \mathbf{Z}_q = \mathbf{I}$. From this, it is clear that this problem admits a global minimum x , and this x satisfies

$$\begin{aligned} & \nabla \left(\text{Tr}(\mathbf{A}^\top \mathbf{A}) - 2x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x + (x^\top \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x) \|x\|_2^2 \right) = 0 \\ \implies & (2 - \|x\|_2^2) \mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x = \|\mathbf{A} \mathbf{Z}_q x\|_2^2 x. \end{aligned} \quad (2)$$

Left multiplying x^\top yields

$$2(1 - \|x\|_2^2) \|\mathbf{A} \mathbf{Z}_q x\|_2^2 = 0.$$

If $\mathbf{A} \mathbf{Z}_q x = 0$, then $\|\mathbf{A} - \mathbf{A} w w^\top\|_F^2 = \|\mathbf{A}\|_F^2$. Otherwise, $\|x\|_2 = 1$. Plugging into equation (2),

$$\mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q x = \|\mathbf{A} \mathbf{Z}_q x\|_2^2 x,$$

which means x is a (unit) eigenvector of $\mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q$ with a non-zero eigenvalue. Since \mathbf{Z}_q is orthonormal, $y = \mathbf{Z}_q x$ is also a unit vector, which means yy^\top is a rank-1 projection matrix. Therefore, the problem is equivalent to maximizing

$$\|\mathbf{A} y y^\top\|_F^2 = \|\mathbf{A} y\|_2^2 = \|\mathbf{A} \mathbf{Z}_q x\|_2^2,$$

which is achieved when x is the top eigenvector of $\mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q$. Since \mathbf{Z}_q is full rank and orthonormal, and \mathbf{A} is symmetric, $\mathbf{Z}_q^\top \mathbf{A}^\top \mathbf{A} \mathbf{Z}_q = (\mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q)^2$ and $\mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q$ is symmetric. Therefore, x is also the top eigenvector of $\mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q = \mathbf{T}_q$, i.e. $x = u$.

Next, we show that if $q = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\gamma}}\right)$, then there exists a unit vector y in the span of \mathbf{Z}_q such that $|\langle v_1, y \rangle| \geq 1 - \epsilon$.

With some work, it can be shown that \mathbf{Z}_q is indeed an orthonormal basis of the Krylov subspace $\mathbf{K}_q(\mathbf{A}, z_1)$; for a full proof, see [8]. Therefore, for any polynomial p_q of degree at most q there exists an x such that $\mathbf{Z}_q x = p_q(\mathbf{A})z_1$. Suppose we show that there is a good approximate top eigenvector in the Krylov subspace, that is, there is a polynomial p_q such that $p_q(\mathbf{A})z_1$ is an approximate top eigenvector of \mathbf{A} . Then, from our previous claim about $w = \mathbf{Z}_q u$, the vector w is also an approximate top eigenvector of \mathbf{A} . Note crucially that we only need to show the existence of such a polynomial, do not need to explicitly compute it.

To this end, let $z_1 = \sum_{i=1}^d \mu_i v_i$, where v_1, v_2, \dots, v_d are the eigenvectors of \mathbf{A} corresponding to eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d \geq 0$. Then,

$$p_q(\mathbf{A})z_1 = \left(\sum_{i=1}^d p_q(\lambda_i) v_i v_i^\top \right) \left(\sum_{i=1}^d \mu_i v_i \right) = \sum_{i=1}^d \mu_i p_q(\lambda_i) v_i.$$

The goal is to find p_q such that $p_q(\lambda_1)$ is large, and $p_q(t)$ is small for any $0 \leq t \leq \lambda_2 \leq (1 - \gamma)\lambda_1$. The following lemma (see Lemma 5 in [5]) on polynomial approximations is helpful:

Lemma 1. *Let $\epsilon', \gamma \in (0, 1)$. Then, there exists a polynomial p of degree at most $O\left(\frac{1}{\sqrt{\gamma}} \log \frac{1}{\epsilon'}\right)$ such that $p(1) = 1$ and $|p(t)| \leq \epsilon'$ for all $0 \leq t \leq 1 - \gamma$.*

We will also require the following high probability bound on $|\mu_1|$:

Lemma 2. *Let $g \sim \mathcal{N}(0, \mathbf{I}_{d \times d})$ and $\delta \in (0, 1)$. Then, with probability at least $1 - \delta$,*

$$\left| \frac{\|g\|_2^2}{d} - 1 \right| = O\left(\frac{\log(2/\delta)}{d} + \sqrt{\frac{\log(2/\delta)}{d}}\right).$$

Moreover, with probability at least $1 - \delta$, $|g_1| = \Omega(\delta)$.

For a proof of Lemma 2 and more concentration inequalities for (sub)-Gaussians, see [9]. Note that our assumption that $\delta = \exp(-O(d))$ implies $\frac{\log(2/\delta)}{d} = O(1)$, so for all sufficiently large $\delta = \exp(-O(d))$, $\|g\|_2 \leq 2\sqrt{d}$ with probability at least $1 - \frac{\delta}{2}$. Therefore, the values μ_i satisfies with probability at least $1 - \delta$,

$$|\mu_1| = |\langle z_1, v_1 \rangle| = \left| \frac{\langle g_1, v_1 \rangle}{\|g\|_2} \right| \geq C \frac{\delta}{\sqrt{d}}$$

for a sufficiently small universal constant C and $|\mu_i| \leq 1$ for $i \geq 2$. Here, we used the fact that $|g_1| = |\langle g, e_1 \rangle|$ and $|\langle g, v_1 \rangle|$ are identically distributed.

Let $\epsilon' \leq \frac{C\delta\sqrt{\epsilon}}{d}$ and $p_q(t) := \hat{p}_q\left(\frac{t}{\lambda_1}\right)$ where \hat{p}_q is the polynomial promised by Lemma 1. Then, $p_q(\lambda_1) = 1$ and $|p_q(\lambda_i)| \leq \epsilon'$ for all $2 \leq i \leq d$. Letting $\rho_i := \mu_i p_q(\lambda_i)$,

$$\frac{|\rho_i|}{|\rho_1|} = \frac{|\mu_i| \epsilon'}{|\mu_1|} \leq \frac{\epsilon' \sqrt{d}}{C\delta} \leq \sqrt{\frac{\epsilon}{d}}.$$

This implies that for $q = O\left(\frac{1}{\sqrt{\gamma}} \log \frac{1}{\epsilon'}\right) = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\gamma}}\right)$,

$$\frac{|\langle p_q(\mathbf{A})z_1, v_1 \rangle|^2}{\|p_q(\mathbf{A})z_1\|_2^2} = \frac{\rho_1^2}{\rho_1^2 + \rho_2^2 + \dots + \rho_d^2} \geq \frac{\rho_1^2}{\rho_1^2 + (d-1)\epsilon \rho_1^2} \geq \frac{1}{1 + \epsilon} \geq 1 - \epsilon.$$

It follows that the polynomial $\frac{p_q(\mathbf{A})z_1}{\|p_q(\mathbf{A})z_1\|_2}$, and therefore the vector w , satisfies equation (1). ■

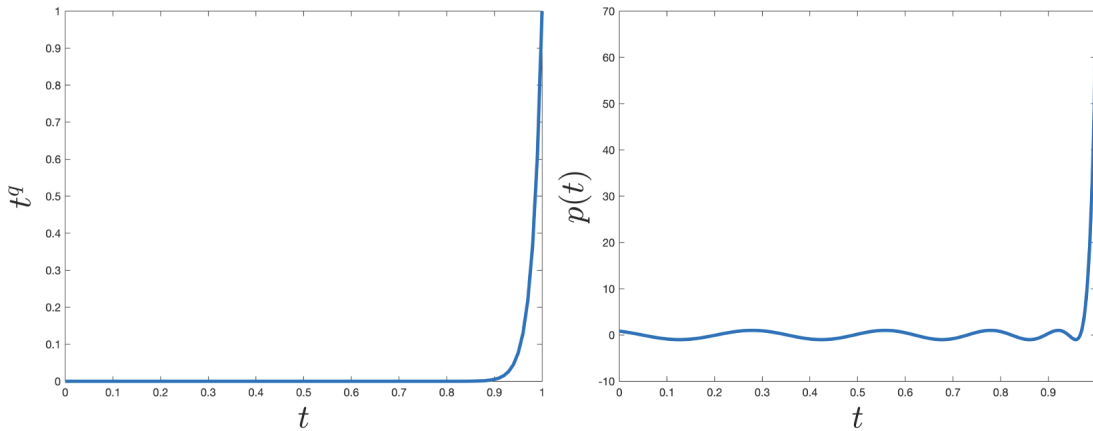


Figure 1: Comparison between t^q and $p(t)$, where p is the (unscaled) polynomial guaranteed by 1.

Remark 2. Note the $\sqrt{\gamma}$ improvement in the the runtime of the Lanczos iteration over the Power iteration 1. The main step achieving this improvement is the polynomial p_q used in 1. The Power Method applies the same technique using the polynomial $f(t) = t^q$. However, the dependence of q on γ is worse:

$$(1 - \gamma)^t \leq \epsilon' \implies t = \Omega\left(\frac{1}{\gamma} \log \frac{1}{\epsilon'}\right).$$

It turns out that t^q can be approximated with a polynomial of degree roughly \sqrt{q} . See [4, 6, 2] for more details.

3.1 Block Krylov Methods

In the previous lecture, we saw the Block Power Method for computing the top k singular vectors of \mathbf{A} . We can similarly extend the Lanczos algorithm to the *Block Lanczos Algorithm*, which leads to a similar quadratic improvement in the number of iterations the algorithm.

Algorithm 3: Block Lanczos Algorithm

Data: $\mathbf{A} \in \mathbb{R}^{n \times d}, q \in \mathbb{N}, k \in \mathbb{N}$

- 1 Choose a random Gaussian matrix $\mathbf{S} \in \mathbb{R}^{d \times k}$
 - 2 $\mathbf{K} \leftarrow [\mathbf{S}, \mathbf{A}\mathbf{S}, \dots, \mathbf{A}^{q-1}\mathbf{S}]$
 - 3 $\mathbf{Z} \leftarrow \text{orth}(\mathbf{K})$, an orthogonal basis of \mathbf{K}
 - 4 $\mathbf{T} \leftarrow \mathbf{Z}_q^\top \mathbf{A} \mathbf{Z}_q$
 - 5 $\tilde{\mathbf{U}}_k \leftarrow$ top k eigenvectors of \mathbf{T}
 - 6 **return** $\mathbf{Z}_q \tilde{\mathbf{U}}_k$
-

Theorem 4 (Block Lanczos Algorithm). *Let \mathbf{V}_k be the top k eigenspace of $\mathbf{A} \in \mathbb{S}_{\geq 0}^{d \times d}$. Let $\epsilon, \delta \in (0, 1)$ with $\delta = e^{-O(d)}$. If the Block Lanczos's algorithm (Algorithm 3) is initialized with $q = O\left(\frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\epsilon}}\right)$, then the output $\mathbf{Z} := \mathbf{Z}_q \tilde{\mathbf{U}}_k$ satisfies*

$$\|\mathbf{A} - \mathbf{A} \mathbf{Z} \mathbf{Z}^\top\|_F^2 \leq (1 + \epsilon) \|\mathbf{A} - \mathbf{A} \mathbf{V}_k \mathbf{V}_k^\top\|_F^2 \quad (3)$$

with probability at least $1 - \delta$. Moreover, the algorithm takes time $O\left(\text{nnz}(\mathbf{A}) \frac{\log(d/\epsilon) + \log(1/\delta)}{\sqrt{\gamma}} k\right)$.

4 Linear System Solvers

We pick up from our motivation of Krylov subspaces to solve linear systems. Given a nonsingular matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ and a vector $b \in \mathbb{R}^d$, solve the system

$$\mathbf{A}x = b. \quad (4)$$

When the matrix \mathbf{A} does not enjoy a particular structure, iterative methods are one of the most popular ways to finding an approximate solution. The idea is to solve for x via updates of the form

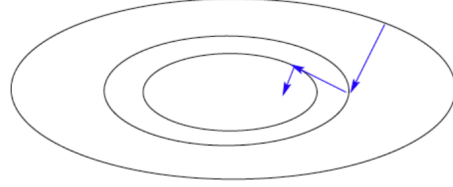
$$x_{\ell+1} \leftarrow x_\ell + \alpha r_\ell$$

for some scalar α and direction vectors r_ℓ which depend on the initial vector x_0 . One such Krylov subspace method is MINRES (Minimum Residual Method). The idea is to pick x_ℓ to be the vector in the (shifted) Krylov subspace $x_0 + \text{Span}(r_0, \mathbf{A}r_0, \dots, \mathbf{A}^{\ell-1}r_0)$ (where $r_0 = b - \mathbf{A}x_0$) which minimizes $\|b - \mathbf{A}x_\ell\|_2$:

$$x_\ell \leftarrow \arg \min_{x \in x_0 + \text{Span}(r_0, \mathbf{A}r_0, \dots, \mathbf{A}^{\ell-1}r_0)} \|b - \mathbf{A}x\|_2$$

This is equivalent [7, 3] to moving along the direction of steepest descent:

$$\begin{aligned} r_\ell &\leftarrow b - \mathbf{A}x_\ell \\ \alpha &\leftarrow \frac{\langle r_\ell, r_\ell \rangle}{\langle \mathbf{A}r_\ell, r_\ell \rangle} \\ x_{\ell+1} &\leftarrow x_\ell + \alpha r_\ell \end{aligned}$$



Here, we highlight the Lanczos method for solving linear systems for symmetric matrices (a similar method exists for non-symmetric matrices via Arnoldi's iterations), which can be viewed as a repeated projection onto the Krylov subspace $\mathbf{K}_q(\mathbf{A}, b - \mathbf{A}x_0)$ and equivalent to the steepest descent method for solving linear systems.

Algorithm 4: Lanczos Algorithm for Linear Systems

Data: $\mathbf{A} \in \mathbb{R}^{n \times d}, b \in \mathbb{R}^d, x_0 \in \mathbb{R}^d, q \in \mathbb{N}$

- 1 $r_0 \leftarrow b - \mathbf{A}x_0, \beta_1 \leftarrow \|r_0\|, r_0 \leftarrow r_0/\beta_1$
 - 2 **for** $\ell = 1, 2, \dots, q$ **do**
 - 3 $y_\ell \leftarrow \mathbf{A}z_\ell - \beta_\ell z_{\ell-1}$
 - 4 $\alpha_\ell \leftarrow \langle y_\ell, z_\ell \rangle$
 - 5 $y_\ell \leftarrow y_\ell - \alpha_\ell z_\ell$
 - 6 $\beta_{\ell+1} \leftarrow \|y_\ell\|_2$. If $\beta_{\ell+1} = 0$ then exit the loop.
 - 7 $z_{\ell+1} \leftarrow \frac{y_\ell}{\beta_{\ell+1}}$
 - 8 $\mathbf{Z}_q \leftarrow [z_1 \ z_2 \ \dots \ z_q], \mathbf{T}_q \leftarrow \text{tridiag}(\beta_j, \alpha_j, \beta_{j+1})$
 - 9 $x_q \leftarrow x_0 + \mathbf{Z}_q \mathbf{T}_q^{-1}(\beta_1 e_1)$
 - 10 **return** x_q
-

4.1 Conjugate Gradient Method

The conjugate gradient (CG) method is a popular variant of the Lanczos algorithm for linear system, when the matrix \mathbf{A} is positive semidefinite. In exact arithmetic, the Lanczos algorithm and the conjugate gradient method are identical.

If the matrix \mathbf{A} is well-conditioned with condition number κ , then the CG method guarantees:

$$\|x^* - x_q\|_{\mathbf{A}} \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^q \|x^* - x_0\|_{\mathbf{A}}$$

Algorithm 5: Conjugate Gradient Method

Data: $\mathbf{A} \in \mathbb{S}_{\geq \mathbf{0}}^{d \times d}$, $b \in \mathbb{R}^d$, $x_0 \in \mathbb{R}^d$, $q \in \mathbb{N}$

```
1  $r_0 \leftarrow b - \mathbf{A}x_0$ ,  $p_0 \leftarrow r_0$ 
2 while the algorithm has not converged, do
3    $\alpha_\ell = \langle r_\ell, r_\ell \rangle / \langle \mathbf{A}p_\ell, p_\ell \rangle$ 
4    $x_{\ell+1} \leftarrow x_\ell + \alpha_\ell p_\ell$ 
5    $r_{\ell+1} \leftarrow r_\ell - \alpha_\ell \mathbf{A}p_\ell$ 
6    $\beta_\ell \leftarrow \langle r_{\ell+1}, r_{\ell+1} \rangle / \langle r_\ell, r_\ell \rangle$ 
7    $p_{\ell+1} \leftarrow r_{\ell+1} + \beta_\ell p_\ell$ 
```

References

- [1] Ed S Coakley and Vladimir Rokhlin. A fast divide-and-conquer algorithm for computing the spectra of real symmetric tridiagonal matrices. *Applied and Computational Harmonic Analysis*, 34(3):379–414, 2013.
- [2] Petros Drineas, Ilse CF Ipsen, Eugenia-Maria Kontopoulou, and Malik Magdon-Ismail. Structural convergence results for approximation of dominant subspaces from block krylov spaces. *SIAM Journal on Matrix Analysis and Applications*, 39(2):567–586, 2018.
- [3] Carl T Kelley. *Iterative methods for optimization*. SIAM, 1999.
- [4] John C Mason and David C Handscomb. *Chebyshev polynomials*. Chapman and Hall/CRC, 2002.
- [5] Cameron Musco and Christopher Musco. Randomized block krylov methods for stronger and faster approximate singular value decomposition. *Advances in neural information processing systems*, 28, 2015.
- [6] Christopher Musco. Singular value decomposition and krylov subspace methods. <https://www.chrismusco.com/amlds2023/notes/lecture11.html>.
- [7] Yousef Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.
- [8] Lloyd N Trefethen and David Bau. *Numerical linear algebra*. SIAM, 2022.
- [9] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.