

**Iterative solvers for linear systems**  
**APPM 4600 Project**  
UNIVERSITY OF COLORADO BOULDER  
DEPARTMENT OF APPLIED MATHEMATICS

## 1 Project Summary

The cost of solving a linear system of the form  $\mathbf{Ax} = \mathbf{b}$  is the step that often dominates the cost of in most algorithms. This is because if  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the cost of inverting  $\mathbf{A}$  directly cost  $O(n^3)$ . It is for this reason that in practice, it is very common to use an iterative solver to construct an approximate solution of the linear system. Roughly speaking, iterative solvers create a sequence of approximations to the solution. The goal is for the sequence to converge to the solution. In class, you learned about Conjugate Gradient which is an iterative solver but only works for matrices with special structure. In this project, you will explore several more widely applicable iterative solvers. You will learn about there potential and limitations.

## 2 Project Background

This project is concerned with approximating the solution to

$$\mathbf{Ax} = \mathbf{b}$$

where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  and  $\mathbf{b} \in \mathbb{R}^n$  via an iterative solver. These iterative solvers only require the ability to apply the matrix  $\mathbf{A}$  to a vector. Typically, these solvers are coupled with fast matrix multiplication methods that either exploit sparsity or structure in the matrix to reduce the cost of applying the matrix to a vector from  $O(n^2)$  with dense linear algebra to  $O(n)$ .

### 2.1 Richardson iteration

The first iterative solver you will investigate is called *Richard iteration*. Roughly speaking the iteration technique recast the problem as a root finding task.

For this section, assume that  $\mathbf{A}$  is positive definite or positive semidefinite. Let  $\mathbf{x}^*$  denote the exact solution to the linear system. This means that  $\mathbf{Ax}^* = \mathbf{b}$ . The Richard iteration is given by the following:

$$\begin{aligned} \text{given} \quad & \mathbf{x}_0 \\ \mathbf{x}_{k+1} &= (\mathbf{I} - \alpha \mathbf{A}) \mathbf{x}_k + \alpha \mathbf{b} \end{aligned}$$

where  $\alpha$  is some user defined constant.

### 2.1.1 Questions to Investigate

1. Show that  $(\mathbf{I} - \alpha \mathbf{A}) \mathbf{x}^* + \alpha \mathbf{b} = \mathbf{x}^*$ .
2. Define the absolute error at step  $k+1$  to be  $\|\mathbf{x}^* - \mathbf{x}_{k+1}\|$ . What conditions are required in order for the method to converge? With these conditions how should you choose  $\alpha$ ? Are there any conditions on  $\mathbf{x}_0$ ?

*Hint: Start by considering  $\mathbf{A}$  is also symmetric and remember our fixed point iteration theory from class.*

3. Test the performance of the method for a collection of matrices that satisfy the specific conditions listed for the method. For which matrices is the method faster than inverting directly? Does this depend on size, sparsity pattern, etc.?

*Hint: you will want to learn how to use the Python fast matrix multiplication toolbox.*

## 2.2 GMRES: Generalized Minimum Residuals

The second iterative method you will consider is GMRES. This method is probably the most widely used iterative solver as it can be utilized for general matrices. While it is relatively easy to gain intuition on how to make the method converge quickly, conducting the actual analysis explaining how the method converges, etc. is still an open problem. This section will introduce you to the method and give you some exercises to gain intuition on the behavior of the method.

### 2.2.1 Krylov subspaces

The *Krylov subspaces*  $\mathcal{K}_m$  generated by  $\mathbf{A}$  and  $\mathbf{b}$  are defined as follows

$$\mathcal{K}_m = \{\mathbf{b}, \mathbf{A}\mathbf{b}, \dots, \mathbf{A}^{m-1}\mathbf{b}\}.$$

In other words, it is the space spanned by the vectors involving powers of  $\mathbf{A}$  being applied to  $\mathbf{b}$ .

A matrix  $\mathbf{A}$  can be written in Hessenberg form via the following orthogonal similarity transformation

$$\mathbf{A} = \mathbf{Q}\mathbf{H}\mathbf{Q}^*$$

where  $\mathbf{Q}$  is orthogonal and  $\mathbf{H}$  is a *Hessenberg* matrix. Let  $\mathbf{Q}_m$  denote the  $n \times m$  matrix whose columns are the first  $m$  columns of  $\mathbf{Q}$ . Let  $\mathbf{q}_j$  denote the  $j^{\text{th}}$  column of  $\mathbf{Q}_m$ . These vectors form a basis for  $\mathcal{K}_m$ .

Let  $\mathbf{K}_m$  denote the matrix

$$\mathbf{K}_m = [\mathbf{b} | \mathbf{A}\mathbf{b} | \mathbf{A}^2\mathbf{b} | \dots | \mathbf{A}^{m-1}\mathbf{b}]$$

Then the QR factorization of  $\mathbf{K}_m$  is given by

$$\mathbf{K}_m = \mathbf{Q}_m \mathbf{R}_m$$

where  $\mathbf{Q}_m$  is the matrix defined above.

### 2.2.2 Now to GMRES

The goal of GMRES is pretty simple. Essentially the goal is to construct  $\mathbf{x}_k$  so that the norm of the residual vector  $\mathbf{r}_k$  is minimized where

$$\mathbf{r}_k = \mathbf{b} - \mathbf{A}\mathbf{x}_k.$$

This means that the linear system solve has been changed to a least squares problem. One option for finding the minimizing vector would be to find  $\mathbf{c} \in \mathbb{R}^n$  such that

$$\|\mathbf{A}\mathbf{K}_k\mathbf{c} - \mathbf{b}\|$$

is minimized. Unfortunately, this is both unstable and expensive.

Instead the vector  $\mathbf{x}_k$  is expressed as  $\mathbf{x}_k = \mathbf{Q}_k\mathbf{y}$  where  $\mathbf{y}$  is the minimizer of

$$\|\mathbf{A}\mathbf{Q}_k\mathbf{y} - \mathbf{b}\|.$$

The algorithm is given as follows (take from [?]) :

1. *Start* with an initial guess  $\mathbf{x}_0$  and compute the residual  $\mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0$  and the normalized residual vector  $\mathbf{v}_1 = \frac{\mathbf{r}_0}{\|\mathbf{r}_0\|}$ .
2. *Iterate:* for  $j = 1, 2, \dots, k, \dots$  until "satisfied" do:
  - $h_{i,j} = \langle \mathbf{A}\mathbf{v}_j, \mathbf{v}_j \rangle$  for  $i = 1, 2, \dots, j$  (Make Hessenberg matrix)
  - $\hat{\mathbf{v}}_{j+1} = \mathbf{A}\mathbf{v}_j - \sum_{i=1}^j h_{i,j}\mathbf{v}_i$  (project away from space already captured)
  - $h_{j+1,j} = \|\hat{\mathbf{v}}_{j+1}\|$
  - $\mathbf{v}_{j+1} = \frac{\hat{\mathbf{v}}_{j+1}}{h_{j+1,j}}$
3. *Form the approximation*

$$\mathbf{x}_k = \mathbf{x}_0 + \mathbf{V}_k\mathbf{y}_k$$

where  $\mathbf{V}_k$  is the matrix with  $j^{\text{th}}$  column vector  $\mathbf{v}_j$  and  $\mathbf{y}_k$  is the minimizer of  $\|\beta\mathbf{e}_1 - \mathbf{H}_k\mathbf{y}\|$ . Here  $\beta = \|\mathbf{r}_0\|$ ,  $\mathbf{e}_1$  denotes the first canonical vector and  $\mathbf{H}_k$  denotes the Hessenberg matrix with entries  $h_{i,j}$  defined by the algorithm.

### 2.2.3 Questions to investigate

1. Write a code that finds the minimizer of

$$\|\beta\mathbf{e}_1 - \mathbf{H}_k\mathbf{y}\|.$$

Here  $\beta = \|\mathbf{r}_0\|$ ,  $\mathbf{e}_1$  denotes the first canonical vector and  $\mathbf{H}_k$  denotes the Hessenberg matrix with entries  $h_{i,j}$  defined by the algorithm.

2. Come up with some ways to set a stopping criterion for GMRES.  
*Convergence error, value of  $h_{j,j}$*  Provide motivation for each of your potential stopping conditions.

3. Write your own GMRES code. Which of your potential stopping conditions works best and why?
4. Compare the number of iterations, accuracy, etc. of your GMRES with the GMRES from a software package. What are differences and similarities? Which would you use in practice?
5. It is pretty easy to construct a matrix where you have control over the condition number and the spectrum. Let  $\mathbf{R}$  denote an  $n \times n$  matrix with random entries. This matrix is full rank. Let  $\mathbf{Q}$  denote the orthogonal matrix that results from the QR factorization of  $\mathbf{R}$ . (You could just use  $\mathbf{R}$  but I like having control over my eigenvectors.) Let  $\mathbf{D}$  denote a diagonal matrix where you get to choose the diagonal entries  $\{d_1, \dots, d_n\}$  (ordered largest to smallest). Then  $\mathbf{A} = \mathbf{Q}^* \mathbf{D} \mathbf{Q}$  is a matrix that you know all the spectral information about.

Fix  $n = 2000$ . Fixing the desired accuracy of GMRES at  $1e - 10$ , solve  $\mathbf{A}\mathbf{x} = \mathbf{b}$  where  $\mathbf{b}$  is vector with random entries and  $\mathbf{A}$

- (a) has  $n$  distinct well separated spectrum.
- (b) is full rank but only has 3 distinct eigenvalues. 1 distinct eigenvalue.
- (c) is full rank and all eigenvalues live in a ball of radius  $1e - 5$  centered at 1.
- (d) is ill-conditioned. i.e. condition number  $1e20$ .
- (e) has one zero eigenvalue but  $\mathbf{b} \in \text{Range}(\mathbf{A})$ .

What have you learned from these experiments? How does the number of required iteration change for each problem?

Does this give you ideas of when GMRES performs well? Does it give you an idea of what you should be looking for in a preconditioner?

### 3 Software Expectations

It is expected that all of the code used in this project is your own except when the directions tell you otherwise or you have gotten approval to use a software package.

### 4 Independent Directions

Possible next steps for the project include, but are not limited to:

1. Investigating and reporting on different iterative methods for approximating eigenvalues.
2. Investigating and reporting on additional iterative solvers such as the Biconjugate Gradient Stabilized Method or multigrid.

3. Present how these methods are used in practice. Note: you will have to get professor approval to ensure your plan is mathematical enough.

## References

- [1] , G. Opfer and G. Schober Richardson's iteration for nonsymmetric matrices *Linear Algebra and its Applications*, 58: 343–361, 1984.
- [2] Y. Saad and M.H. Schultz GMRES: A Generalized Minimal Residual Algorithm for Solving Nonsymmetric Linear Systems *SIAM Journal on Scientific and Statistical Computing*, 7(3): 856–869, 1986.