Mathematical Institute

University of Oxford

# Implementing a Matrix Class to Solve Linear Systems and Find Eigenvalues

# Contents

# 1    Introduction

We implement a class of matrices and subclass of vectors in C++, overloading opera-
tors where necessary to allow for robust functionality. Upon implementing a definite
linear system solver in Gaussian elimination, and indefinite solver in GMRES, we
show their correctness by solving Poisson's equation for different forcing functions us-
ing the finite element method. We then consider eigenvalue problems, implementing
the famous $QR$ algorithm. We test our eigenvalue solver on an ODE based eigenvalue
problem, again using the finite element method.

# 2    Matrices in C++

We begin my introducing the class of matrices that form the basis of this report. Inside
our header file, **Matrix.h** we have defined our **protected**[1] and **public** variables and
methods, key definitions can be seen in the following code.

```
7   class Matrix
8   {
9   protected:
10          double* matrixVal = nullptr;
11          int rows;
12          int columns;
13  public:
14          // Used to allocate space
15          void createMatrix(int rows, int columns);
16          // Constructors
17          Matrix();
18          Matrix(int r, int c);
19          Matrix(const Matrix& m);
20          Matrix(int size) :Matrix(size, size) {};
21          // Destructor
22          ~Matrix();
```

Lines **11-12** show the expected **protected** variables for a class of matrices. Perhaps
less conventionally, on line **10** we have defined a pointer to the first entry of a single
array of doubles, rather than to the first entry of an array of pointers which each
point to an array of doubles. In this construction we store a $m \times n$ matrix in memory
as a single array of doubles of length $mn$, as opposed to the latter method which

---

[1]We note that some of our variables are defined **protected** instead of **private** as they must be
accessible by our inheriting vector sub-class

would store the same matrix in memory as $m$ distinct arrays of doubles of size $n$, a seemingly pointless difference however it is justified by the speed increase, which we shall demonstrate after introducing the constructors and destructor. Line **15** is the method called by the constructors defined on lines **17-19** to dynamically allocate memory, and line **20** allows for a single size parameter to be specified which allocates memory for a square matrix[2]. Line **22** defines our destructor to free the memory. The following code shows the **createMatrix** method located in our **Matrix.cpp** source file.

```cpp
void Matrix::createMatrix(int rows, int columns)
{
    // Delete the contents of any non-empty matrix
        if (matrixVal != nullptr)
        {
                delete[] matrixVal;
        }
        // Allocate memory for single array format
        matrixVal = new double[rows * columns];
        this->rows = rows;
        this->columns = columns;
        if (rows == 0 || columns == 0)
        {
            // Assign an empty matrix NULL value
                matrixVal = NULL;
        }
        else
        {
                for (int i = 0; i < rows * columns; i++)
                {
                    // Assign the matrix with zeros
                        matrixVal[i] = 0;
                }
        }
}
```

The key line that demonstrates our method of storing the matrix is line **13**, where we simply allocate an array of doubles of the required size. The default constructor **Matrix()**[3] then simply calls **createMatrix(0,0)** and the basic overload of this,

---

[2]this simply calls the constructor that allocates memory for generic matrices with rows = columns = size, and so need not be defined in **Matrix.cpp**.

[3]We make use of this format of default constructor so we can define a matrix whose size we may not want to specify.

**Matrix(r,c)** calls **createMatrix(r,c)**. The copy constructor is similar, creating a matrix of the same size and simply assigning the value at each index.

## 2.1 Cache Friendly Memory Allocation and Locality

As stated already, we allocate memory using a pointer to the first entry in a single array of doubles of length $m \times n$, instead of a pointer to an array of $m$ pointers which each point to the first entry in an array of doubles of length $n$. This is a matter of data locality[4], i.e. the locality of data entries which are likely to be accessed together. In the case where the matrices are stored as pointers to arrays of pointers, the data is not contiguous in memory since there is no reason why the secondary pointers should point to entries in any structured manner. In contrast to this, the method we employ points to the first entry in single array which is by definition contiguous in memory - a more cache friendly approach. Consequently, in any iteration over the matrix, the CPU avoids cache misses. In other words the CPU avoids having to wait for data that is not currently cached due to the lack of structure in the underlying data. It also allows the CPU to prefetch data[1], which helps reduce the natural bottleneck of transferring data from memory to the CPU. We can indeed demonstrate this, and to do so we implement a timer function. The code for this function can be found in appendix C.3, and it is aptly called by executing **tic();** prior to the code and **toc();** afterwards. We record the time taken for both methods in allocating memory and assigning values. For a matrix of size $8000 \times 8000$, and averaged over 30 executions, we find that the average time taken by our method is 0.1862 seconds, compared to 0.3288 seconds in the opposing method when timed on the same machine - a significant speed increase. A small downside to this method is that for a **Matrix** $A$, we are required to index using **[i+A.rows*j]** where we would otherwise use **[i][j]**. We refer the reader to this fact when reading the code displayed in this report.

## 2.2 Operator Overloading

We have overloaded a series of operators to allow for robust functionality. All of the basic mathematical binary $(+, -, /, *)$ and unary $-$ operators have been overloaded to allow for computations between types **double**, **Matrix** and **Vector**. The equals $(=)$ operator has been overloaded to allow for the assigning of **Matrix** types (and also **Vector** types )to each other. Similarly, the parentheses $((\cdot \, , \, \cdot))$ operator has

---

[4]also referred to as spatial locality or memory locality

been overloaded to allow for the indexing of a **Matrix** or **Vector** , which we allow to start from 1. Finally the $<<$ operator has been overloaded to allow for the streaming of these types. A full list of these, as well as all implemented function, can be seen in our header file - see appendix C.1.

# 3    Solving $Ax = b$

An ample application of our matrix class is to solve some linear system $Ax = b$ that appears in a mathematical problem. In this section we discuss the algorithms implemented to solve such a problem, we consider these algorithms as described by Trefethen and Bau in [7], by Layton and Sussman in [3] and for GMRES in particular, by Saad[5] in [5] .

## 3.1    Gaussian Elimination

We first consider the implementation of a direct method to solving linear systems in Gaussian elimination (with partial pivoting to ensure stability). The full code for this can be found in appendix C.4. The implementation can be broken down into three main steps. Given an $m \times m$ matrix $A$ and a $m \times 1$ Vector $b$, we seek the solution $x$ to $Ax = b$. The first step is to form the augmented matrix $\tilde{A} = [\,A\,|\,b\,]$, where we have simply appended the vector $b$ as an additional column on $A$. Once this has been formed, we perform suitable row operations to transform the augmented matrix into row reduced form. Finally we employ back substitution to solve for the vector $x$. In practice, the augmented matrix is formed by appending the $m \times 1$ vector $b$ to the right of the $m \times m$ matrix $A$. This is a convenient step especially for the programming implementation since, in essence, it allows us to map the row transformations that we apply to $A$ directly onto $b$ so we need not store these in a separate matrix. We complete this task by introducing two new functions to the **Matrix** class, **addCols(int c)** and **setColumn(int col, Vector& v)**, these will also prove useful in the implementation of GMRES. We simply call the former to add a single column to (a copy of) $A$ and the latter to append the vector $b$ onto this new column. Once the augmented matrix has been computed we begin our loop through its columns. The first objective in our loop is to find the largest (absolute) value below the diagonal of the current column, we then use our implemented **swapRows** function, which takes in the indices of two rows, to swap this with the subdiagonal

---

[5]Saad is one of the co-publishers for the paper that introduced GMRES

entry (and thus its row). We call this new entry the pivot, and it is used to eliminate all of the elements below it. If the value of the largest pivot is 0, then the current column need not be reduced so the column index is increased and we move onto the next column. As Trefethen and Bau note, this method of partial pivoting introduces backwards stability to Gaussian elimination and is thus superior to the standard, non-pivoting algorithm [7]. Once we have the correct (non zero) pivot, the remainder of the loop simply performs row operations to reduce the other elements in the column to zero, thus changing the entries in the corresponding rows in the remainder of the matrix. After looping through all columns we arrive at a row reduced form, for example the final form when $A$ is $4 \times 4$ may be

$$\left[\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & a_{14} & b_1 \\ & a_{22} & a_{23} & a_{24} & b_2 \\ & & a_{33} & a_{34} & b_3 \\ & & & a_{44} & b_4 \end{array}\right],$$

where we note that all non labelled elements are definitely zero, however those labelled could take any value (a zero element on the diagonal would constitute a singular matrix, in which case we throw an exception). From this point we compute the entries of the solution vector $x$ by performing back substitution [3]. Using the above matrix as an example, back substitution would start with assigning $x_4 = b_4/a_{44}$, and then the row above would yield the equation $a_{33}x_3 + a_{34}x_4 = b_3$ and one could plug in the calculated $x_4$ and solve for $x_3$. This process is repeated in a backwards loop (starting at the final row and finishing at the first), appending values found to the **Vector** $x$ at each stage. Once the loop has completed we have found our solution so we return its value.

## 3.2 GMRES

In contrast to direct algorithms, which aim to factorise the matrix $A$ into a product of more simple matrices, thereafter immediately arriving at the exact solution , iterative algorithms rely on gradually refining the solution by minimising some form of error [7]. Krylov subspace methods are a subcategory of iterative algorithms that rely on approximating the solution to $Ax = b$ in some space defined by polynomials of $A$ times the vector $b$. In particular we have that the $nth$ Krylov subspace is defined by

$$\mathcal{K}_n = \langle b, Ab, \cdots, A^{n-1}b \rangle, \tag{1}$$

where $\langle \cdot \rangle$ denotes the linear span. We consider the implementation of the Generalised Minimal Residual Method [5] (GMRES), which is such an iterative method.

At iteration $k$ of GMRES, we seek an an approximate solution $x_k \in \mathcal{K}_k$ to $Ax = b$ that minimises the residual $r_k = ||r_0 - Ax_k||_2$, that is $x_k = \text{argmin}_{\hat{x} \in \mathcal{K}_k}||r_0 - A\hat{x}||_2$, where $r_0 = b - Ax_0$, and $x_0$ is an initial guess to the solution, often take as a vector of zeros. We implement a for loop to search for such solutions. We begin at iterate $k = 1$, where $k$ denotes the loop for finding the minimising vector in the $k$th Krylov subspace. The termination condition is $k <$ **maxits**, where **maxits** is a user specified maximum allowed iterations of the algorithm. In each loop we want to search for a optimal solution in the $k$th Krylov subspace and so need a basis to check in. To this end we employ Arnoldi iteration to construct an orthonormal basis of the subspace ([5],[7]). The convenience in this method is that it generates a single vector on the current iteration using the vector generated from the previous iteration, where the vectors are mutually orthonormal - repeating this process generates, at iteration $k$, an orthogonal matrix $Q_{k+1}$ whose columns are the orthonormal basis vector of $\mathcal{K}_{k+1}$. Along with this, Arnoldi iteration also generates an upper Hessenberg matrix $\tilde{H}_k$ at iteration $k$ via the relation

$$AQ_k = Q_{k+1}\tilde{H}_k. \tag{2}$$

This Hessenberg matrix will be crucial for the remainder of the algorithm. Prior to our loop we set $Q$ to be a **Matrix**[6] of size $m \times 1$, where $m$ is the side length of the square matrix $A$. The first orthonormal basis vector of $Q$ is set to be $r_0/||r_0||$ and we proceed to enter the main loop of GMRES. Within this main loop, we begin our Arnoldi iteration loop for the current GMRES iterate $k$. This follows the standard Arnoldi iteration algorithm [7], as listed in Appendix A.1 .

```
1        // Arnoldi starts here ————————————————————————
2        // q_previous stores the previous column of q
3                q_previous = getColumn(Q, k − 1);
4                Q.addCols(1);
5                // At this point H is (k+1)x(k) and Q is mxk
6        // Initially set next column as A*(previous column)
7                q = A * q_previous;
8        // Initialise next column in the Hessenberg matrix
9                Vector h(k + 1);
10       // Main Arnoldi Loop ————————————————————————
```

---

[6]An important note here is that we set $Q$ to be a **Matrix** of size $m \times 1$ instead of a **Vector** of length $m$, this is because throughout the loop we append our newest othonormal basis vector as an additional column, thereby changing the shape.

```
11                      for ( int  i = 0;  i < k;  i++)
12                      {
13          // This  follows  the  standard  Arnoldi  algorithm
14                          qi = getColumn (Q, i );
15                          qt = transpose (q );
16                          double  hi = ( qt * qi ). matrixVal [ 0 ];
17                          h. matrixVal [ i ] = hi ;
18                          q = q − hi * qi ;
19                      }
20          // ——————————————————————————————
21          // Assign  the  final  value  of  new H  column
22                  h. matrixVal [ k ] = norm ( q );
23          // Final  value  of  the  newest  column  of Q
24                  q = q * (1.0 / norm ( q ));
25          // add  the  new  basis  vector  to  the  matrix
26                  Q. setColumn ( k + 1,  q );  // ———————————> Q_{k+1}
27          // add  new  column  to  Hessenberg  matrix
28                  H. setColumn ( k,  h );       // ———————————>  H_{k}
29          // Arnoldi  ends  here ——————————————————————————
```

Line **7** simply initialises the newest orthonormal basis vector **q** as the product of $A$ and the previous basis vector, which is taken from the orthonormal matrix $Q$ where these are stored. The vector which shall be appended to this iterations' Hessenberg matrix is initialsied on line **9** . Lines **11-19** perform the Arnoldi iteration to compute the newest vectors to be appended to $Q_{k+1}$ and $\tilde{H}_k$, and the code on lines **26** and **28** does exactly that. Upon creating these matrices, $Q_{k+1}$ and $\tilde{H}_k$, we can re-wrtite the objective of GMRES in this loop. By setting $x_k$, the required approximate solution in this iteration, to a linear combination of the basis vectors[7] $Q_k$ (i.e. the the first $k$ columns of $Q_{k+1}$), $x_k = x_0 + Q_k y$, our goal is to find the minimiser $y$ to the $k$th error, or residual; $||r_k||_2 = ||Ax_0 + AQ_k y - b||_2 = ||AQ_k y - r_0||_2 = ||Q_{k+1}\tilde{H}_k y - r_0||_2$. Where we have used the recurrence relation 2. Since we initialise the first column of $Q$, to be $r_0/||r_0||_2$, it follows that we can write $r_0 = ||r_0||_2 Q_{k+1} e_1$, where $e_1$ is the first standard basis vector of $\mathbb{R}^m$ (1 in the first entry and 0s in the remaining $m - 1$ entries). Thus we want to find the minimiser $y$ to $||r_k||_2 = ||Q_{k+1}\tilde{H}_k y - ||r_0||_2 Q_{k+1} e_1||_2 = ||\tilde{H}_k y - ||r_0||_2 e_1||_2$, by the unitary invariance [7] of the Euclidean norm and the orthogonality of $Q_{k+1}$. All of this derivation to say - we can solve for $x_k$ by solving for coefficients of its basis

---

[7]Recall that the goal of GMRES is simply that - to find the optimal solution in the Krylov subspace, and these vectors span it so naturally we construct $x$ as a linear combination of them and find the optimal coefficients, given by $y$. The addition of $x_0$ is included since we are calculating residuals relative to some initial guess $x_0$.

expansion, $y$, and this boils down to solving a least squares problem involving known quantities - the $(k + 1) \times k$ Hessenberg matrix computed from Arnoldi iteration, $\tilde{H}_k$, and the term easily calculated from the initial residual, $||r_0||_2 e_1$. In practice we solve the least squares problem by computing the $QR$ factorisation of $\tilde{H}_k$, a task whose cost is very cheap due to the structure of an upper Hessenberg matrix being almost upper triangular already. For brevity and sake of repetition, we omit the details of the $QR$ factorisation process here, but note that it will be covered in section 4 where we discuss the $QR$ algorithm for finding eigenvalues. The important part for GMRES is that we can essentially compute what are known as Givens rotation matrices[5] for $\tilde{H}_k$, that upon left multiplication, knock off the subdiagonal entries turn by turn, thus forming an upper triangular matrix. Of course $\tilde{H}_k$ is not square; it is $(k + 1) \times k$, and so the resulting matrix would be a square, upper triangular matrix with an extra row of zeros which we can easily remove. We note that this process is done iteratively, and throughout the iterations we track the total product of all Givens rotation matrices as these will constitute the relevant transformation from Hessenberg to upper triangular on the termination of the loop. For brevity we omit the code for the decomposition of $\tilde{H}_k$, but note that it can be found between lines **100-193** in appendix C.5. The main GMRES loop terminates if the current error, calculated by a product of the previous errors and the current Givens rotation matrix, is lower than some user specified tolerence, **tol**, or if, as stated previously the iterations surpass the maximum allowed. Say we terminate at iteration $k = N$. Upon termination we simply take our elongated, $(N + 1) \times N$ upper triangular matrix, **UT** in the code, and remove the final row, consisting of zeros, naming the new **Matrix U**, for upper triangular. We then take $Q_N$ by removing the final column of the currently stored $Q_{N+1}$. We use Gaussian elimination to solve $Uy = \beta$, where we have tracked residuals in the **Vector** $\beta$. This problem is very cheap since $U$ is upper triangular and thus the cost is only in the necessary back substitution. Upon finding $y$ we easily calculate the approximation $x_N = x_0 + Q_N y$. The implementation of computing the final value is fairly trivial and can be found between lines **237-273** in appendix C.5.

### 3.2.1 Convergence of GMRES

We can check our implementation of GMRES by investigating the number of iterations required for the solution of $Ax = b$ for increasingly well conditioned matrices $A$. We apply GMRES to an $M \times M$ matrix of random entries uniformly distributed between 0 and 1, using our implemented **rand(int n)** function. This takes in an integer **n**

and returns a square **Matrix** of side length **n** filled with such random entries. For the sake of this experiment, we temporarily modify our **gmres** function to return the **Vector** of errors, $||r_k||/||b||$ over a specified number of iterations. We then compute the errors for our random matrix $A$, and for $A + nI$, with $n = 10, 20, \cdots, 100$. We consider Theorem 35.2 by Trefethen and Bau, as stated in [7]. We first note that the $k$th residual can be written $||r_k|| = ||p_k(A)b||$, where $p_k \in P_k$. Here $P_k$ is the set of polynomials of degree (at most) $k$. Then, the theorem states that at the $k$th iteration of GMRES, $r_k$ satisfies

$$\frac{||r_k||}{||b||} \leq \kappa(V) \inf_{p_k \in P_k} \sup_{\lambda \in \Lambda(A)} |p_k(\lambda)|, \tag{3}$$

where $\Lambda(A)$ is the set of eigenvalues of $A$ and $V$ is the matrix of eigenvectors of $A$. Without going into the technical details of the results of the theorem, it roughly tells us that we expect the convergence to speed up when eigenvalues of the matrix are clustered away from the origin. Trefethen and Bau go on to give examples[7] of random matrices whose conditioning is improved by adding on multiples of the identity (i.e. clustering the eigenvalue spectrum of $A$ away from the origin). We demonstrate similar results, shown in figure 1. We see a drastic decrease in iterations required for the convergence of our **gmres** function. The code used to create this figure can be found in appendix C.5.1.



Figure 1: Convergence results for a random $100 \times 100$ matrix $A$, taking **gmres(A,b,guess,1.0e-10,1000)**, where $b$ is a **Vector** of ones, **guess** is a **Vector** of zeros, and the final two entries are the tolerance and maximum iterations respectively. The arrow shows the direction of increasing $n$, for values $n = 0, 10, 20, \cdots, 100$

.

11

## 3.3 Application: Finite Element Solution to the Poisson Equation

We test our linear solvers on the Poisson equation with homogeneous Dirichlet boundary conditions and yet unspecified forcing function $f$. We follow the derivation[8] of the discretization by [6]. We want to find $u : \Omega \mapsto \mathbb{R}$ that satisfies

$$-\Delta u = f \text{ in } \Omega, \tag{4}$$

$$u = 0 \text{ on } \Gamma \tag{5}$$

numerically. Multiplying by a test function $v \in V$ (where $V$ is some function space which we shall choose to match the requirements of the variational formulation) and integrating gives

$$-\int_\Omega v \Delta u \, dx = \int_\Omega v f \, dx. \tag{6}$$

Applying the divergence theorem and noticing that $\nabla \cdot (v \nabla u) = v \Delta u + \nabla u \cdot \nabla v$ yields

$$-\int_\Omega v \Delta u \, dx = \int_\Omega \nabla u \cdot \nabla v \, dx - \int_\Gamma v \nabla u \cdot n \, ds, \tag{7}$$

where $n$ is the unit outward normal on $\Omega$ to $\Gamma$ and $ds$ is a line element on $\Gamma$. Now, by choosing $V = H_0^1(\Omega)$, the Sobolev space that permits once weakly differentiable functions (so we can work with $\nabla u$ and $\nabla v$) that satisfy the boundary condition (5), $v$ vanishes on $\Gamma$ and so the final term in (7) vanishes. We now have the variational (weak[9]) form of our problem; find $u \in V$ such that

$$\int_\Omega \nabla u \cdot \nabla v \, dx = \int_\Omega v f \, dx \tag{8}$$

for all $v \in V$. Though a weakened formulation, we are still trying to find a solution in a infinite dimensional function space $V$, thus we are seemingly still not in sight of a numerical solution, let alone a linear system to solve. The key is in the next step. We search for solutions in a finite dimensional subspcae $V_h$ of $V$. Such a task is described by the Galerkin approximation[6] to our problem; find $u_h \in V_h$ such that

$$\int_{\Omega_h} \nabla u_h \cdot \nabla v_h \, dx = \int_{\Omega_h} v_h f \, dx \tag{9}$$

for all $v_h \in V_h$. We look for solutions $u_h$ with general $v_h$ in a basis of $V_h$. Recalling the definition of $V = H_0^1(\Omega)$, we require once (weakly) differentiable functions that vanish

---

[8]The derivation will also be relevent when solving the eigenvalue problem in the next section.

[9]We also refer to this as the weak formulation, since we have weakened the requirements on $u$ as a solution to our problem.

on the boundary. We define a mesh of a domain to be a series of non intersecting cells whose union is the entire space. Considering $V_h$ as such a mesh, we define a nodal basis $V_h = \text{span}\{\phi_1, \cdots, \phi_N\}$ of functions $\phi_i$ which vanish both on the boundary and on the intersection points between cells of the mesh, called the degrees of freedom. Writing $u_h$ and $v_h$ as a linear combination of basis functions, $u_h = \sum_{i=1}^N U_i \phi_i$ and $v_h = \sum_{j=1}^N V_j \phi_j$. After plugging these into (9) and rearranging, we arrive at

$$\sum_{j=1}^N V_j \sum_{i=1}^N U_j \int_{\Omega_h} \nabla \phi_i \nabla \phi_j dx = \sum_{j=1}^N V_j \int_{\Omega_h} \phi_j f dx. \tag{10}$$

We can write this as $AU = F$ where

$$A_{ij} = \int_{\Omega_h} \nabla \phi_i \nabla \phi_j dx, \quad \text{and} \quad F_j = \int_{\Omega_h} \phi_j f dx, \tag{11}$$

and $U$ is the vector of coefficients that will be our discrete solution. Choosing the one-dimensional domain $\Omega = [0,1]$ (and thus the boundary conditions $u(0) = u(1) = 0$), we can take the uniform mesh as a linear spacing of $[0,1]$, with spacing $h$. With this formulation we can consider the domain as a mesh defined by degrees of freedom at interior points with spacing $h$, and we can consider the discretized solution space $V_h$ as having dimension equal to the number of interior mesh points, $N = 1/h - 1$, where the numerator is simply the domain length, and we take one less to account for the extra exterior node included in the quotient. We omit the discussion regarding the rigorous formulation of the basis functions and the master element mapping for brevity however, essentially, we evaluate our integrals over the master element $[0,1]$ instead of over each local finite element defined by degrees of freedom. In doing this we also calculate a local matrix for $A$ and local vector for $F$ defined by the mapping between the local element and the master element. This is done via the transformation $x \mapsto \xi = (x - x_L)/(x_R - x_L)$, where the local element defined by degrees of freedom $x_L$ and $x_R$ is mapped to the master element defined by $[0,1]$. We note that in this case our basis functions become simply $\phi_1(\xi) = 1 - \xi$ and $\phi_2(\xi) = \xi$ for the one dimensional case. Applying the master element transformation, the matrix assembly requires that we compute the local matrix

$$\tilde{A} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} = \begin{bmatrix} \int_0^1 h^{-1} \phi_1' \phi_1' d\xi & \int_0^1 h^{-1} \phi_1' \phi_2' d\xi \\ \int_0^1 h^{-1} \phi_1' \phi_2' d\xi & \int_0^1 h^{-1} \phi_2' \phi_2' d\xi \end{bmatrix} \tag{12}$$

for each cell and append these $2 \times 2$ matrices on the main diagonal of $A$. Then the global matrix $A$ is assembled as

$$
A = \begin{bmatrix}
\color{blue}{a_{11}} & \color{red}{a_{12}} & & & & & \\
a_{21} & \mathbf{a_D} & a_{12} & & & & \\
 & a_{21} & \mathbf{a_D} & a_{12} & & & \\
 & & a_{21} & \mathbf{a_D} & \ddots & & \\
 & & & \ddots & \ddots & & \\
 & & & & & \mathbf{a_D} & a_{12} \\
 & & & & & \color{red}{a_{21}} & \color{red}{a_{22}}
\end{bmatrix},
\tag{13}
$$

where we have defined $\mathbf{a_D} = a_{22} + a_{11}$. We note that in contrast to $A$, the computation of $F$ is dependent on a forcing function $f$, to which we also must apply the mapping. We have that $f(x) = f(x_L + \xi(x_R - x_L))$, thus we must define $F$ in a more general sense. We have the local vector

$$
\tilde{F}_i = \begin{bmatrix} F_i \\ F_{i+1} \end{bmatrix} = \begin{bmatrix} \int_0^1 h\phi_1 f(x_i + \xi(x_{i+1} - x_i)d\xi \\ \int_0^1 h\phi_2 f(x_i + \xi(x_{i+1} - x_i)d\xi \end{bmatrix},
\tag{14}
$$

where $x_i$ and $x_{i+1}$ are the values of the degrees of freedom $i$ and $j$. We then have the global vector assembled as

$$
F = \begin{bmatrix}
\color{red}{F_1} \\
\mathbf{F_2 + F_3} \\
\mathbf{F_4 + F_5} \\
\vdots \\
\mathbf{F_N - 2 + F_N - 1} \\
\color{red}{F_N}
\end{bmatrix}.
\tag{15}
$$

We note that the entries in red shall be set to 0, and the entries in blue shall be set to 1 so as to satisfy the boundary conditions. From here we can implement the problem in to our **main()** method using the built up **Matrix** class. For convenience, we implement a **linspace** function, which takes as parameters two **double** types denoting the initial values and final values respectively, and an **int** type which denotes the number of points in our spacing. This function returns a **Vector** object which stores the linear spacing. The following code then generates the required variables for 100 mesh points.

```
1  int meshpts = 100;
2  Vector mesh = linspace(0.0,1.0,meshpts);
```

14

```
 3    Vector  dofs = linspace(1,meshpts,meshpts);
 4    double h = 1.0 / ((double) meshpts − 1.0); // spacing
 5    // Initialise A matrix (stiffness matrix) − square of size meshpts
 6    Matrix A(meshpts);
 7    // Initialise F vector (forcing vector) − size meshpts
 8    Vector F(meshpts);
 9    // Initialise DOFS
10    int dof1 , dof2 ;
```

We next assemble our matrix $A$ and vector $F$. By noticing the symmetry of $\tilde{A}$, we need only compute $a_{11}, a_{12}$ and $a_{22}$, and these values are constant so do not need to be evaluated in a loop. To calculate integrals over $[0, 1]$ we apply Simpson's rule[10], namely

$$\int_0^1 p(x)dx \approx \frac{1}{6}\left(p(0) + 4p\left(\frac{1}{2}\right) + p(1)\right).  \tag{16}$$

Thus, by defining functions $f, \phi_1, \phi_2, \phi_1', \phi_2',$

```
 1    double f(double x){return −1;}   // Example forcing function f(x) = −1
 2    double phi1(double x){return 1.0 − x;}
 3    double phi2(double x){return x;}
 4    double gradphi1(double x){return −1.0;}
 5    double gradphi2(double x){return 1.0;}
```

we can compute the three distinct elements of $\tilde{A}$ by using 16, a task whose code we omit for brevity though refer the reader to appendix C.6. We then loop through the elements of the mesh, bar the last, and let our degrees of freedom be the current iterate and the subsequent iterate in the loop. These degrees of freedom denote the indices in the global matrix $A$ to which we append our submatrix $\tilde{A}$ (and similarly for $F$). The matrix and vector assembly is then simply written in the loop as the following code.

```
 1         for (int element = 1; element < dofs(meshpts); element++)
 2         {
 3              // dofs give the index in A that we add to
 4              dof1 = element;
 5              dof2 = element + 1;
 6              // evaluate Simpsons rule points on the local cells
 7              x1 = masterToLocal(0.0, mesh(dof1), mesh(dof2));
 8              x2 = masterToLocal(0.5, mesh(dof1), mesh(dof2));
 9              x3 = masterToLocal(1.0, mesh(dof1), mesh(dof2));
```

---

[10]Of course we could easily calculate these integrals by hand and simply assign the values, however Simpson's rule will be necessary in the eigenvalue problem example and thus we introduce it here. It also turns out that in this example, the Simpson's rule approximation is exact

```
10                    // evaluate F elements at correct dofs
11                    // using Simpsons rule
12                    F1 = (double)1.0 / 6.0 * (h * f(x1) * phi1(0.0)) +
13                            +4.0 / 6.0 * (h * f(x2) * phi1(0.5)) +
14                            +1.0 / 6.0 * (h * f(x3) * phi1(1.0));
15                    F2 = (double)1.0 / 6.0 * (h * f(x1) * phi2(0.0)) +
16                            +4.0 / 6.0 * (h * f(x2) * phi2(0.5)) +
17                            +1.0 / 6.0 * (h * f(x3) * phi2(1.0));
18                    // Matrix and vector assembly
19                    A(dof1, dof1) += A11;
20                    A(dof1, dof2) += A12;
21                    A(dof2, dof1) += A12;
22                    A(dof2, dof2) += A22;
23                    F(dof1) += F1;
24                    F(dof2) += F2;
25            }
```

In lines **7-9** we have used our mapping between the master element and the local element. Upon completion of the assembly loop, we assign the boundary conditions by setting the values of $A$ and $F$ to 0 and 1 in the relvent locations (seen in red and blue previously). At this point in the code we have our final matrix $A$ and vector $F$ and thus it is finally time to put our linear solvers to the test. We compute solutions using both Gaussian elimination,

```
1            Vector U_ge(meshpts);
2            U_ge = A / F;
```

and GMRES,

```
1            Vector U_gm(meshpts);
2            Vector guess = 0 * U_gm;
3            int maxits = meshpts;
4            double tol = 1.0e-10;
5            U_gm = gmres(A, F, guess, tol, maxits);
```

where we have manually set the tolerance to be $10^{-10}$ and max iterations to be the size of the matrix (GMRES should converge by this point). We then use **std::ofstream** to stream[11] our results to a **.dat** file which we can import into MATLAB to verify our results. For the forcing function $f(x) = -1$, it is not hard to see that we expect the solution $u(x) = x(x-1)/2$. Sure enough we see the correct results for both linear solvers.

---

[11]This is possible since we have overloaded the $<<$ operator to stream **Matrix** types.
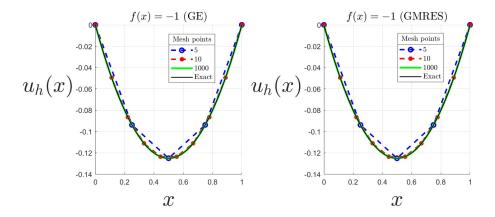
Figure 2: Finite element approximations to the Poisson equation with forcing function $f(x) = -1$. We see solutions for both Gaussian elimination and GMRES being used to solve the resulting linear system, at $5, 10$ and $1000$ mesh points. We also see the exact solution

Now that we have verified our solvers for a simple forcing function, we can see how it holds out on a more complicated example. An easy way to do this is to pick a $u$, solve for $f$ and use that as our forcing function. Choosing $u(x) = 2/3(x^5 - x)$ satisfies the boundary conditions and gives $f = -40/3x^3$. Changing our code to:

```
1  double f(double x){return −40.0 / 3.0 ∗ pow(x,3);}
```

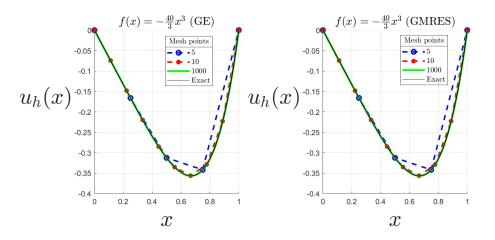and rerunning for mesh sizes $5, 10$ and $1000$ gives us back our solution.



Figure 3: Finite element approximations to the Poisson equation with forcing function $f(x) = -40/3x^3$. We see solutions for both Gaussian elimination and GMRES being used to solve the resulting linear system, at $5, 10$ and $1000$ mesh points. We also see the exact solution $u(x) = 2/3(x^5 - x)$

# 4  Solving $Ax = \lambda x$

Though seemingly similar to $Ax = b$, the solution to the eigenvalue problem $Ax = \lambda x$ requires a different approach. We implement the (pure) $QR$ algorithm [7] which at first glance is strikingly simple - see appendix 2. Essentially, the algorithm repeats finding the $QR$ decomposition of $A$, multiplying $RQ$, finding the $QR$ decomposition of the result and so on. This converges to the Shur form of $A$ [7]. Once this has been computed, the entries on the main diagonal are the eigenvalues. The only difficulty of the implementation of this lies within the step of actually computing the $QR$ factorisation, for clarity we create a separate function, **qr**, to do this. An important note here is the return type. We would like to return both the orthogonal **Matrix** $Q$ and the upper triangular **Matrix** $R$ from the function that computes the decomposition, thus we use the **tuple** return type. After including $\langle$**tuple**$\rangle$ , the declaration in our header file[12] **Matrix.h** reads

```
1            friend std::tuple<Matrix, Matrix> qr(const Matrix& A);
```

and one possible way we could call the function is **std::tie(Q, R) = qr(A)**. The actual return statement in the **qr** functin reads

```
1            return std::make_tuple(Q, R);
```

Onto the implementation of the $QR$ decomposition. We use the Givens rotation matrices to remove non zero entries from $A$ to eventually make it upper triangular [7]. Since Givens rotation matrices are orthogonal, the product and inverse of them is also orthogonal, and thus the process of multiplying by them shall also yield an orthogonal matrix [3]. We initialise $Q = I_m$ and $R = A$ where $A$ is the $m \times m$ matrix we are interested in. The Givens rotation matrix that we use to zero a subdiagonal entry $(i, j)$ in $A$ is given by $\mathrm{G}_{ij}, where$

$$
\underbrace{\begin{bmatrix} I_{i-2} & & & \\ & c & s & \\ & -s & c & \\ & & & I_{m-i} \end{bmatrix}}_{G_{ij}\ (m \times m)} \times \underbrace{\begin{bmatrix} * & \dots & * & \dots & * \\ & & \vdots & & \\ \vdots & \dots & \mathbf{t} & \dots & \vdots \\ \vdots & \dots & \mathbf{b} & \dots & \vdots \\ & & \vdots & & \\ * & \dots & 0 & \dots & * \end{bmatrix}}_{A\ (m \times m)} = \begin{bmatrix} * & \dots & * & \dots & * \\ & & \vdots & & \\ \vdots & \dots & * & \dots & \vdots \\ \vdots & \dots & \mathbf{0} & \dots & \vdots \\ & & \vdots & & \\ * & \dots & 0 & \dots & * \end{bmatrix} \tag{17}
$$

---

[12]we also declare the function outside the scope of the **Matrix** class so as to allow for access in our **main()**.

and $m \geq i \geq j + 1 \geq 2$ since we are dealing with subdiagonal entries. The lower right $c$ of the Givens submatrix is entry $(i, i)$ in the larger matrix[13]. In (17) we have noted the effected rows by the transformation in red. Upon calculating $c$ and $s$ using

$$c = \frac{t}{\sqrt{t^2 + b^2}}, \quad s = \frac{b}{\sqrt{t^2 + b^2}}, \qquad (18)$$

we implement (17) with the following code in our loop through subdiagonal entries in all columns. We note that in each loop, prior to this calculation, we initialise **givens = eye(m)**.

```
1          // assign values to givens matrix
2          givens.matrixVal[i+m*i] = c;
3          givens.matrixVal[i+m*(i-1)] = -s;
4          givens.matrixVal[(i-1)+m*i] = s;
5          givens.matrixVal[(i-1)+m*(i-1)] = c;
6          // apply givens matrix to zero subdiagonal entries
7          R = givens * R;
8          // track the orthogonal vectors
9          Q = Q * transpose(givens);
```

This process looped over all subdiagonal elements will convert $R$ (initialised as $A$) into an upper triangular matrix, and right multiplying $Q$ (initialised as **eye(m)**) by the transpose of the Givens matrix will yield the corresponding orthogonal matrix. This is due to the relations [7]

$$G_{m-1}G_{m-2}\cdots G_1 A = R \qquad (19)$$

$$\implies A = G_1^T \cdots G_{m-2}^T G_{m-1}^T R = QR, \qquad (20)$$

which holds since each Givens matrix is orthogonal, $G^{-1} = G^T$, and thus the product of the inverses (transposes) of Givens matrices is also orthogonal. Over the course of the loop through all subdiagonal elements, line **7** in the code above generates equation (19) and line **9** generates equation (20).

## 4.1   QR algorithm

Now that we can compute the $QR$ decomposition of a square matrix, we can implement the $QR$ algorithm[7] stated in appendix 2. We do this by creating another

---

[13]It is important to note that this matrix seemingly only depends only on the row of the entry we want to zero, however we have implicitly defined the Givens rotation matrix with respect to the entry to be removed and the entry in the column above ($b$ and $t$, for bottom and top) , thus the $2 \times 2$ structure of the submatrix is preserved. If we were considering zeroing with respect to a different entry (not directly adjacent) then our matrix would depend on $j$.

function, **eigMatrix**, to handle the generation of the Schur form of $A$, and a final, MATLAB styled function **eig** which extracts the diagonal elements (eigenvalues) from the Schur form and returns a **Vector** containing them. In practice this implementation is easily done in a loop since we have our **qr** function, and so we can repeatedly find the $QR$ decomposition of the previously calculated product $R \times Q$. We terminate once our Schur form has been found - that is when the matrix is diagonal (for symmetric entries) or upper triangular (for general entries), which we can detect by checking the size of the off diagonals[7]. The **eig**[14] function then simply creates a **Vector** of size equal to that of the rows (or columns) of the matrix and assigns the values from the diagonal entries. We note that in this construction eigenvalues are returned in decreasing order, just as in MATLAB. We can now simply print the eigenvalues of a matrix $A$ with the following code in our **main()**.

```
1  std::cout << "The eigenvalues of " << A << " are " << eig(A) << "\n";
```

## 4.2  Application: Finite Element Solution to an Eigenvalue Problem

We note that the purpose of the $QR$ algorithm is to find eigenvalues, and though some eigenvectors may be generated, we limit our discussions to finding the former. An interesting extension would be to implement a more efficient way to calculate eigenvectors. To verify our eigenvalue calculator, we consider finding the eigenvalues $\lambda_n$ that satisfy the equation[15][2]

$$-u''(x) + c(x)u(x) = \lambda u(x) \tag{21}$$
$$u(a) = u(b) = 0. \tag{22}$$

We use a similar finite element discretisation as when considerring the Poisson equation, as such we omit the derivation of the generalised eigenvalue problem here, but refer the reader instead to appendix B. We end up with the task to find eigenvalues that satisfy

$$AU = \lambda MU. \tag{23}$$

Left multiplying both sides by $M^{-1}$, gives the eigenvalue problem

$$M^{-1}AU = \lambda U. \tag{24}$$

---

[14]**eig** calls **eigMatrix** which repeatedly calls **qr**.

[15]For the case $c = 0$, this is known as the one dimensional Laplace eigenvalue problem [2].

Thus, to find our eigenpairs, we need to assemble the matrices $A$ and $M$, compute the inverse[16] $M^{-1}$, and then find eigenvalues of the product $M^{-1}A$. The boundary conditions are incorporated by setting the entries in the first and last rows of $A$ to zero and the same for $M$ except the first entry of the first row and the final entry of the final row. This has the equivalent effect to enforcing that the product $M^{-1}A$ when right multiplied by $U$, zeros the initial and final values, thereby satisfying the homogeneous Dirichlet boundary conditions. We consider the case $a = 0$ and $b = \pi$. Using the master element transformation as before, we obtain

$$\tilde{a}_{ij} = \int_0^1 \frac{1}{h}\phi_i'\phi_j' + h\phi_1\phi_2 c(x_I + \xi(x_{I+1} - x_I))\,d\xi \qquad (25)$$

where $i = 1, 2$, $j = 1, 2$ and $h = \pi/(N-1)$. Here $x_1$ and $x_2$ denote the $x$ values of the degrees of freedom $I$ and $I + 1$, where the local matrix

$$\tilde{A} = \begin{bmatrix} \tilde{a}_{11} & \tilde{a}_{12} \\ \tilde{a}_{21} & \tilde{a}_{22} \end{bmatrix}$$

is added onto the global matrix at indices $(I, I), (I, I+1), (I+1, I)$ and $(I+1, I+1)$. Clearly $I$ takes values between 1 and $N - 1$, where $N$ is the number of degrees of freedom. The matrix $M$ is defined only in terms of basis functions and therefore all necessary values can be constructed prior to a loop, we have

$$\tilde{M} = \begin{bmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{bmatrix} = \begin{bmatrix} \int_0^1 \phi_1\phi_1 h\,d\xi & \int_0^1 \phi_1\phi_2 h\,d\xi \\ \int_0^1 \phi_2\phi_1 h\,d\xi & \int_0^1 \phi_2\phi_2 h\,d\xi \end{bmatrix}. \qquad (26)$$

We implement the solution to this problem in a very similar manner to the Poisson equation, and the full code can be found in appendix C.8. We note that for the coded solution, it was also convenient to implement an inverse function to calculate $M^{-1}$. We did this using our Gaussian elimination operator and calculating $A/e_i$ where $A$ is the $m \times m$ matrix to be inverted and $e_i$ is the $i$th standard basis function of $\mathbb{R}^m$, iterating $i$ over all rows. By appending the vector obtained from the calculation at each iteration to a new **Matrix** we obtain the inverse. Upon completing the matrix assembly, we simply execute the following code.

```
1    Matrix B = inverse(M) * A;
2    std::cout << "eigenvalues:_" << eig(B) << "\n";
```

---

[16]A glaring issue in the simplification of the generalised eigenvalue problem to the standard eigenvalue problem is the possibility of a singular $M$. We fortunately do not run into this issue here since we choose appropriate functions $c(x)$, however we note that a possible extension would be to implement a solver of generalsied eigenvalue problems, such as the $QZ$ algorithm, see [4].

For $c(x) = 0$, one can see that we require $u(x) = A\sin(\sqrt{\lambda}x) + B\cos(\sqrt{\lambda}x)$. Enforcing $a = 0$ and $b = \pi$ gives $u(x) = sin(\sqrt{\lambda}x)$ with the requirement that $\lambda = 1^2, 2^2, 3^2, \ldots$ ($\lambda_k = k^2$ for $k \in \mathbb{N}$) to satisfy the boundary conditions[17]. Assembling the matrices with these parameters, we find that the above code gives eigenvalues $\lambda_1 = 1.06161, \lambda_2 = 4.98914$ and $\lambda_3 = 12.7081$ with 5 mesh points. This begins to resemble the required result, however we can gain some confidence in our solution when using 50 mesh points; we obtain the following approximations to $\lambda_i$: 1.00034, 4.00548, 9.02778, 16.0882, 25.2225, 36.5188, 50.2117, 66.8098, . . . , 880.766 where we clearly see the correct smaller eigenvalues, and an increasing error for larger eigenvalues.

# 5   Discussion

We have shown that our somewhat basic, from scratch, implementations of these well established numerical linear algebra methods and algorithms have been successful, in that they solve some non trivial problems, on a fairly sized scale. We have demonstrated the robustness of our implementation through our code snippets and have demonstrated a convergence property of our GMRES that gives confidence in our implementation. We have shown that our underlying memory allocation technique conforms to rigorous memory management conventions, in that they aim to minimise cache misses and reduce the bottlneck of the data transfer from memory to the CPU. By no means were the implemented algorithms the optimal choices for all cases, but rather a good choice for most cases. As such a natural extension would be to implement further, more specific solvers that can be chosen for different types of problems in a robust manner.

---

[17]Of course the trivial solution $\lambda = 0$ and thus $u = 0$ would also work.

# References

[1] M. Ananth, S. Vishwas, and M.R. Anala. "Cache Friendly Strategies to Optimize Matrix Multiplication". In: *2017 IEEE 7th International Advance Computing Conference (IACC)*. 2017, pp. 23–27. DOI: `10.1109/IACC.2017.0020`.

[2] Daniele Boffi. "Finite element approximation of eigenvalue problems". In: *Acta Numerica* 19 (2010), 1–120. DOI: `10.1017/S0962492910000012`.

[3] William Layton and Myron Sussman. "Numerical linear algebra". In: *University of Pittsburgh, Pittsburgh* (2014), pp. 28–39.

[4] Cleve B Moler and Gilbert W Stewart. "An algorithm for generalized matrix eigenvalue problems". In: *SIAM Journal on Numerical Analysis* 10.2 (1973), pp. 241–256.

[5] Yousef Saad. *Iterative methods for sparse linear systems*. SIAM, 2003.

[6] "The Finite Element Method for the Poisson Equation". In: *Numerical Methods for Elliptic and Parabolic Partial Differential Equations*. New York, NY: Springer New York, 2003, pp. 46–91. ISBN: 978-0-387-21762-8. DOI: `10.1007/0-387-21762-2_3`. URL: `https://doi.org/10.1007/0-387-21762-2_3`.

[7] Lloyd N Trefethen and David Bau III. *Numerical linear algebra*. Vol. 50. Siam, 1997.

# A    Algorithms

## A.1    Arnoldi Iteration

---

**Algorithm 1:** Arnoldi Iteration [7]

**Data:** Matrix, $A$. Vector, $b$. Initial Guess to $Ax = b$, $x_0$.

**Result:** Orthonormal vectors $q_1, q_2, \cdots$

$r_0 \leftarrow b - Ax_0$

$q_1 \leftarrow r_0 / ||r_0||_2$

**for** $i = 1, 2, 3, \cdots$ **do**

    $v \leftarrow Aq_n$

    **for** $j = 1, \cdots, n$ **do**

        $h_{ji} \leftarrow q_j^T v$

        $v \leftarrow v - h_{ji} q_j$

    **end**

    $h_{i+1,i} \leftarrow ||v||_2$

    $q_{i+1} \leftarrow v / h_{i+1,i}$

**end**

---

## A.2    QR algorithm

---

**Algorithm 2:** Pure QR Algorithm [7]

**Data:** Matrix, $A$.

**Result:** Schur form $\bar{A}$ of $A$

$A_0 \leftarrow A$

**for** $i = 1, 2, 3, \cdots$ **do**

    $Q_i R_i \leftarrow A_{i-1}$                 Compute QR factorisation of $A_{i-1}$

    $A_i \leftarrow R_i Q_i$

**end**

---

# B Derivation of eigenvalue problem

We consider the discretisation of the following problem.

$$-u''(x) + c(x)u(x) = \lambda u(x) \tag{27}$$

$$u(a) = u(b) = 0 \tag{28}$$

where $u_n \in [a, b]$, $c(x) \geq 0$ and $a < b$. We follow a very similar discretisation process as we did for the Poisson equation in section 3.3. Multiplying by a test function in a currently unknown function space $v \in V$, and integrating over our domain gives

$$\int_a^b -u''v + cuv \, dx = \lambda \int_a^b uv \, dx. \tag{29}$$

Integrating by parts and enforcing $V = H_0^1([a, b])$ then gives our weak form of the eigenvalue problem; find $u \in V = H_0^1([a, b])$ such that

$$\int_a^b u'v' + cuv \, dx = \lambda \int_a^b uv \, dx \tag{30}$$

for all $v \in V = H_0^1([a, b])$. Recall that this holds since the Sobolev space $H_0^1(\Omega)$ enforces that functions vanish on the boundary of $\Omega$ and the space permits functions whose first weak derivative exists, thus (30) is a well defined problem. We then, as before, search for a solution in the finite dimensional subspace $V_h \subset V$. Taking $u_h$ and $v_h$ as linear combinations of basis vectors $\phi_1, \phi_2, \cdots$. We have the problem to find the Galerkin approximation $u_h$ to the eigenvalue problem by finding coefficients $U_i$ of the basis expansion of $u_h$ such that

$$\sum_{j=1}^N V_j \sum_{i=1}^N U_i \int_a^b \phi_i'\phi_j' + \phi_i\phi_j c \, dx = \lambda \sum_{j=1}^N V_j \sum_{i=1}^N U_i \int_a^b \phi_i\phi_j \, dx \tag{31}$$

for all coefficients $V_j$ of the basis expansion of $v_h$. Since this must hold for all coefficients $V_j$, it is equivalent to finding the vector of coefficients $U$ that satisfies the generalised eigenvalue problem

$$AU = \lambda MU. \tag{32}$$

where

$$A_{ij} = \int_a^b \phi_i'\phi_j' + \phi_i\phi_j c \, dx \tag{33}$$

and

$$M_{ij} = \int_a^b \phi_i\phi_j \, dx. \tag{34}$$

25

# C  C++ Code

## C.1  Matrix.h

```cpp
1  #ifndef MATRIXDEF
2  #define MATRIXDEF
3
4  #include <cmath>
5  #include "Exception.h"
6  #include <tuple>
7  class Vector; //forward declaration
8  class Matrix
9  {
10 protected:
11         double* matrixVal = nullptr;
12         int rows;
13         int columns;
14 public:
15
16         // Used to allocate space
17         void createMatrix(int rows, int columns);
18         // Constructors
19         Matrix(); // Default constructor
20         Matrix(int r, int c);// Constructor given rows and columns
21         Matrix(int size) :Matrix(size, size) {}; // Given size
22         Matrix(const Matrix& m); // Copy constructor
23         // Destructor
24         ~Matrix();
25
26         // Operators
27         Matrix& operator=(const Matrix &m);      // equality
28         double& operator()(int i, int j);        // indexing
29         friend Matrix operator+(const Matrix& m1, const Matrix& m2);
30         friend Matrix operator-(const Matrix& m1, const Matrix& m2);
31         friend Matrix operator*(const Matrix& m1, const Matrix& m2);
32         friend Matrix operator*(const double& a, const Matrix& m);
33         friend Matrix operator*(const Matrix& m, const double& a);
34         friend Matrix operator/(const Matrix& m, const double& a);
35         friend Matrix operator-(const Matrix& m);
36         friend Vector operator*(const Matrix& m, const Vector& v);
37         friend Vector operator/(const Matrix& A, const Vector& b);
38         // printing
```

```
39        friend std::ostream& operator<<
40        (std::ostream& output, const Matrix& m); // printing
41        // Functions
42        friend Matrix eye(int n);
43        friend Vector getColumn(Matrix A, int c);
44        friend Matrix getRow(Matrix A, int c);
45        friend Matrix transpose(const Matrix& m);
46        friend Vector gmres(const Matrix& A, const Vector& b,
47            const Vector& x0, double tol, int maxits);
48        friend std::tuple<Matrix, Matrix> qr(const Matrix& A);
49        friend Matrix eigMatrix(const Matrix& A);
50        friend Vector eig(const Matrix& A);
51        friend Matrix inverse(const Matrix& A);
52        friend Matrix hessenberg(const Matrix& A);
53        friend Vector convertMatrixToVector(const Matrix& A);
54        friend Matrix rand(int n);
55        Matrix& swap_rows(int a, int b);
56        Matrix& setColumn(int col, Vector& v);
57        Matrix& addRows(int r);
58        Matrix& addCols(int c);
59
60 };
61 // Accesible functions
62 Matrix rand(int n);
63 Matrix hessenberg(const Matrix& A);
64 Matrix inverse(const Matrix& A);
65 Matrix eigMatrix(const Matrix& A);
66 std::tuple<Matrix, Matrix> qr(const Matrix& A);
67 Vector gmres(const Matrix& A, const Vector& b,
68     const Vector& x0, double tol = 1.0e-10, int maxits = 10000);
69 Matrix eye(int n);
70 Vector eig(const Matrix& A);
71 Vector getColumn(Matrix A, int c);
72 #endif
```

## C.2 Vector.h

```
1 #include "Matrix.h"
2
3 class Vector : public Matrix
4 {
5 public:
```

```
6
7            // Constructors
8            Vector();
9            Vector(int size);
10           Vector(const Vector& v);
11
12           // Operators
13           Vector& operator=(const Vector& v);       // equality
14           double& operator()(int i);                            // indexing
15           friend Vector operator-(const Vector& m1, const Vector& m2);
16           friend Vector operator+(const Vector& m1, const Vector& m2);
17           friend Vector operator*(const Vector& m, const double& a);
18           friend Vector operator*(const double& a, const Vector& m);
19
20           friend double norm(Vector v, int p);
21           friend Vector linspace(double a, double b, int n);
22           friend double length(Vector v);
23
24 };
25 Vector linspace(double a, double b, int n);
26 double norm(Vector v, int p = 2);
27 double length(Vector v);
```

## C.3   Memory Allocation Testing

### C.3.1   Timer.h

```
1  #ifndef TIMERDEF
2  #define TIMERDEF
3
4  #include <chrono>
5
6  static std::chrono::
7  time_point<std::chrono::high_resolution_clock> startTime;
8
9  void tic();
10 void toc();
11
12 #endif // !TIMERDEF
```

### C.3.2   Timer.cpp

```
1  include "Timer.h"
2  #include <stdio.h>
3  void tic()
4  {
5          startTime = std::chrono::high_resolution_clock::now();
6  }
7
8  void toc()
9  {
10         float dt = std::chrono::
11         duration_cast<std::chrono::microseconds>
12         (std::chrono::high_resolution_clock::now()−startTime).count()
13         / 1000000.f;
14         printf("Elapsed time is %f seconds.\n", dt);
15 }
```

### C.3.3   Main

```
1  #include "Timer.h"
2  int main()
3  {
4
5          // we test the allocation of memory
6          // and assigning of entries for the two methods
7          int rows = 8000;
8          int columns = 8000;
9          //  single array storage
10         tic();
11         double* A = new double[rows * columns];
12         for (int i = 0; i < rows * columns; i++)
13         {
14                 A[i] = 0;
15         }
16         delete[] A;
17         toc();
18         // multi array storage
19         tic();
20         double** B = new double*[rows];
21         for (int i = 0; i < rows; i++)
22         {
23                 B[i] = new double[columns];
24         }
25         for (int i = 0; i < rows; i++)
```

29

```
26              {
27                      for ( int  j  =  0;  j  <  columns ;  j++)
28                      {
29                              B[ i ] [ j ]  =  0;
30                      }
31              }
32              for ( int  i  =  0;  i  <  rows ;  i++)
33              {
34                      delete [ ]  B[ i ] ;
35              }
36              delete [ ]  B;
37              toc ( ) ;
38
39
40              return  −1;
41  }
```

## C.4   Gaussian Elimination

```
1  Vector operator /( const  Matrix& A,  const  Vector& b)
2  {
3          int  m = A. rows ;
4          int  n = m +  1;  //  Columns  of  Augmented  matrix
5          //  We  only  accept  square  matrix  input
6          if  (m != A. columns )
7          {
8                  throw  Exception (
9                          ”Bad␣input” ,
10                         ”A/b␣ requires ␣a␣sqaure ␣matrix␣input” );
11         }
12         //  Form  augmented  matrix
13         Vector  bcopy(b );
14         Matrix  A␣aug(A);
15         A␣aug . addCols ( 1 );
16         A␣aug . setColumn (m +  1,  bcopy );
17
18         //  solution  vector
19         Vector  x(m);
20
21         //  Initialise  variables
22         double  largest␣pivot  =  0;
23         int  largest␣pivot␣index  =  0;
```

```
24          double q;           // variable to store a quotient in the loop
25
26          int h = 0; // row index      (0 ---> m - 1)
27          int k = 0; // column index (0 ---> n - 1 = m)
28
29
30          // loop through the augmented matrix
31          while ((h < m) && (k < n))
32          {
33                  // clear previos pivot values
34                  largest_pivot = 0;
35                  largest_pivot_index = 0;
36                  // findng pivot in column k -
37                  // loop through all subdiagonal entries
38                  for (int i = h; i < m; i++)
39                  {
40                          // check if entry > largest pivot
41                          if(abs(A_aug.matrixVal[i+k*m])
42                                      >abs(largest_pivot))
43                          {
44                                  // if entry > largest pivot,
45                                  // set largest pivot = entry
46                                  largest_pivot=A_aug.matrixVal[i+k*m];
47                                  largest_pivot_index=i;
48
49                          }
50
51                  }
52                  // if largest pivot is 0 we can move
53                  // onto the next column
54                  if (largest_pivot == 0)
55                  {
56                          k += 1;
57                  }
58                  // otherwise we bring the pivot to the sub diagonal
59                  // and eliminate lower entries
60                  else
61                  {
62                          A_aug.swap_rows(h+1,largest_pivot_index+1);
63                          for (int i = h + 1; i < m; i++)
64                          {
65                                  // find the ratio of each
66                                  // value to the pivot
```

```
67                              q = A_aug.matrixVal[i+k*m]/
68                                  A_aug.matrixVal[h+k*m];
69                              // elimnate entries below the pivot
70                              A_aug.matrixVal[i+k*m]=0;
71                              for (int j = k + 1; j < n; j++)
72                              {
73                                      // perform corresponding
74                                      // operation on the rest
75                                      // of the row (-=)
76                                      A_aug.matrixVal[i+j*m] -=
77                                      A_aug.matrixVal[h+j*m]*q;
78                              }
79                          }
80                          // move to the next column, repeat the process
81                          h += 1;
82                          k += 1;
83                  }
84
85          }
86
87          // throw exception if matrix is singular
88          for (int i=0;i<m;i++)
89          {
90              if(A_aug.matrixVal[i+i*m] == 0)
91              {
92                  throw Exception("singular_matrix",
93                  "A/b_cannot_solve_when_A_is_singular");
94              }
95          }
96
97          // Now we find x using back substitution
98          double sum = 0;
99
100         // final value of x is simply the ratio of
101         // the final two non zero elements in the row
102         // reduced augmented matrix
103         x.matrixVal[m - 1] = A_aug.matrixVal[m-1+(n-1)*m]/
104                             A_aug.matrixVal[m-1+(n-2)*m];
105         // we loop through the rest of the rows, each time
106         // gaining another value in x
107         for (int i = m - 2; i >= 0; i--)
108         {
109                 sum = 0;
```

```
110                  for (int j = i + 1; j <= m - 1; j++)
111                  {
112                          // sum the known values of x * coefficients.
113                          // -= since these will be subtracted from
114                          // the final column value
115                          sum -= A_aug.matrixVal[i + j * m]
116                                  * x.matrixVal[j];
117                  }
118              // next x value is given by adding the sum of
119              // -(row elements * known x values) and the final
120              // coefficient
121              x.matrixVal[i] = (A_aug.matrixVal[i+m*(n-1)]+sum)
122                              / (A_aug.matrixVal[i+i*m]);
123          }
124          // return the solution
125          return x;
126 }
```

## C.5   GMRES

```
1  Vector gmres(const Matrix& A, const Vector& b, const Vector& x0,
2  double tol, int maxits)
3  {
4
5          int m = A.rows;
6          if (m != A.columns)
7          {
8                  throw Exception(
9                  "Bad_input", "GMRES_requires_a_sqaure_matrix_input");
10         }
11         // Residual given initial guess
12         Vector r = b - A * x0;
13         // We initialise Q as a matrix instead of a
14         // vector since we shall add columns
15         Matrix Q(m, 1);
16         Vector q = r * (1.0 / norm(r)); // first entry of Q
17         Q.setColumn(1, q); // q now added as first column of Q
18
19
20         // Initialise the Hessenberg matrix obtained from Arnoldi
21         Matrix H(2, 1); // starts as a 2x1
22
```

```
23          /// Givens rotations values
24          Matrix GivensTotal = eye(2);
25          Matrix GivensCurrent;
26
27          // variables used for Givens rotations
28          Matrix Hcopy(H);
29          double c;
30          double s;
31          double top;
32          double bot;
33          double div;
34          Matrix UT(2);
35          Vector errors(maxits);
36          errors.matrixVal[0] = norm(r);
37          double error;
38          double bnorm = norm(b);
39
40          // we will often require a temporary matrix
41          Matrix temp;
42
43          // variables for Arnoldi
44          Matrix qt;
45          Vector qi;
46          Vector q_previous;
47
48
49          // for gmres demonstration
50          //Vector residuals(maxits);
51          //int itsreq = 1;
52
53          int k = 1; // errors iterations.
54          // we define this outside the scope  of
55          // the loop so we can access it afterwards
56          for (k; k < maxits; k++)
57          {
58
59                  if (k > 1)
60                  {
61                          // on first iteration:
62                          // Q is mx1 and H is 2x1,
63                          // on further iterations we
64                          // increase the size of H to
65                          // retain Upper Hessenberg form
```

```
66                          H. addRows ( 1 ) ;
67                          H. addCols ( 1 ) ;
68                      }
69
70
71              // Arnoldi starts here ————————————————————————
72              // q_previous stores the previous column of q
73              q_previous = getColumn(Q, k − 1);
74              Q. addCols ( 1 ) ;
75              // At this point H is (k+1)x(k) and Q is mxk
76              // Initially set next column as A∗(previous column)
77              q = A ∗ q_previous ;
78              // Initialise next column in the Hessenberg matrix
79              Vector h(k + 1);
80              // Main Arnoldi Loop ————————————————————————
81              for ( int i = 0; i < k; i++)
82              {
83                      // This follows the standard Arnoldi algorithm
84                      qi = getColumn(Q, i );
85                      qt = transpose (q );
86                      double hi = ( qt ∗ qi ). matrixVal [ 0 ] ;
87                      h. matrixVal [ i ] = hi ;
88                      q = q − hi ∗ qi ;
89              }
90              // ————————————————————————————————————————
91              // Assign the final value of new H column
92              h. matrixVal [ k ] = norm( q ) ;
93              // Final value of the newest column of Q
94              q = q ∗ (1.0 / norm( q ) ) ;
95              // add the new basis vector to the matrix
96              Q. setColumn( k + 1, q ); // ————————————————> Q_{k+1}
97              // add new column to Hessenberg matrix
98              H. setColumn( k , h ) ;      // ————————————————>  H_{k}
99
100             // Arnoldi ends here ————————————————————————
101
102
103
104             // Uncommenting the below code gives a check
105             // for orthonormality between basis vectors :
106             /∗
107             // q1 and q1 are columns of Q
108             Matrix q1 = getRow( transpose (Q), k );
```

```
109                    Vector q2 = getColumn(Q, k+1);
110                    std::cout << "Q is " << Q << " and " << q1 <<
111                    " x " << q2 << " = " << q1 * q2 << "\n";
112                    */
113
114                    // We now solve the minimisation problem ————————————
115
116
117                    // Notes on implementation:
118                    // Need to compute the Givens rotation for
119                    //the input H, this depends on the final
120                    // two non zero entries of column k...
121                    // ... once we have multiplied the current
122                    // H by the total of the previous Givens rotations.
123                    // This will give a matrix...
124                    // ... which is upper triangular with an extra
125                     //row of zeros
126                    // Need to store the value of the previous (total)
127                    // Givens as this is required for the update of
128                    // current givens
129                    // Inital value of total givens is 2x2 identity
130                    // as initial H is 2x1.
131                    // - need to change H and keep a copy
132                    // - want the current givens to be kxk
133
134                    // Initialise the givens matrix that will zero
135                    // sub diagonal entries of the newest column of H
136                    // top left corner of givens is identity.
137                    // Bottom right 2x2 is the interesting part
138                    GivensCurrent = eye(k - 1);
139                    GivensCurrent.addCols(2);
140                    GivensCurrent.addRows(2);
141
142                    // we will change the lower right 2x2 submatrix
143                    // in current givens matrix to hold values
144                    // [ c s ; -s c]...
145                    // ... where c and s are dependent on the final
146                    // two non zero entries of the kth column of
147                    // the current upper triangular form.
148
149                    // We want to store the H of this iteration
150                    // as it is used for the next Arnoldi step,
151                    // therefore make a copy that we change instead
```

```
152                     Hcopy = GivensTotal * H;
153

154                // Need to multiply this by the current
155                // Givens matrix to obtain the current
156                // upper triangular matrix
157                // So we need to form the givens matrix.
158                // This depends on the element on the  ,...
159                //...    diagonal and the one below it.
160                // Call these top and bot
161

162                // Equivalent to QR factorising but
163                // since it is Hessenberg initially we
164                // need to zero much fewer values
165

166                // calculate givens matrix constants for this iteration
167                top = Hcopy.matrixVal[(k - 1) + Hcopy.rows * (k - 1)];
168                bot = Hcopy.matrixVal[(k)+Hcopy.rows * (k - 1)];
169                div = pow(pow(top, 2) + pow(bot, 2), 0.5);
170                c = top / div;
171                s = bot / div;
172

173

174

175                // now we have the values of the 2x2 we can add
176                // these to current givens to give our full
177                // transformation matrix
178

179                GivensCurrent.matrixVal[(k-1)+GivensCurrent.rows*(k-1)]
180                =c;
181                GivensCurrent.matrixVal[(k)+GivensCurrent.rows*(k-1)]
182                =-s;
183                GivensCurrent.matrixVal[(k-1)+GivensCurrent.rows*(k)]
184                =s;
185                GivensCurrent.matrixVal[(k)+GivensCurrent.rows*(k)]
186                =c;
187

188                // we track the composition of all givens rotations
189                temp = GivensCurrent * GivensTotal;
190                GivensTotal = temp;
191

192           // current upper triangular matrix
193                UT = GivensTotal * H;
194
```

```
195                        // append newest residual based on previous
196                        // residual and the givens transformation
197                        errors.matrixVal[k] = errors.matrixVal[k - 1] * (-s);
198                        errors.matrixVal[k - 1] = errors.matrixVal[k - 1] * (c);
199                        error = fabs(errors.matrixVal[k] / bnorm);
200                        // for gmres demo
201                        //residuals.matrixVal[k] = error;
202
203
204                        // check termination condition
205                        if (error < tol)
206                        {
207                                std::cout << "GMRES_CONVERGED_IN_"
208                                << k << "_ITERATIONS_WITH_RESIDUAL_"
209                                << error << "\n";
210
211                                //for gmres demo
212                                //itsreq = k;
213                                //
214                                // exit loop
215                                k = maxits;
216                        }
217                        else
218                        {
219                                // Increase the size of the Givens total
220                                // so it is ready for the next loop
221                                GivensTotal.addCols(1);
222                                GivensTotal.addRows(1);
223                                // New givens total = | |————————————| 0 |
224                                //                    | |old givens total| ...|
225                                //                    | |----------------| 0 |
226                                //                    | 0      ...      0   1 |
227                                GivensTotal.matrixVal
228                                        [(k+1)+GivensTotal.rows*(k+1)]=1;
229
230                        }
231                        // Uncomment this to print the current iteration and error
232                        //std::cout << " Current iteration: " << k
233                        //<< " Error : " << error << "\n";
234
235            }
236            // need to remove final row of temp * H to give UT matrix
237            Matrix U(UT.rows - 1, UT.columns);
```

```
238            Vector beta(U.rows);
239            Matrix Qf(Q.rows, Q.columns - 1);
240            // remove final column of Q (since we terminate with
241            // Q_{k+1} but only the upper triangular matrix
242            // gained from the QR decomposition of H_{k})
243            for (int i = 0; i < Qf.rows; i++)
244            {
245                    for (int j = 0; j < Qf.columns; j++)
246                    {
247                            Qf.matrixVal[i+Qf.rows*j]=
248                                    Q.matrixVal[i+Q.rows*j];
249                    }
250            }
251
252            // Extract the  upper triangular matrix from UT
253            // since it has an extra row of zeros
254            // simultaneously extract all but the last rows
255            // of the vector of residuals
256            for (int i = 0; i < U.rows; i++)
257            {
258                    //yy.matrixVal[i] = y.matrixVal[i];
259                    beta.matrixVal[i] = errors.matrixVal[i];
260                    for (int j = 0; j < U.columns; j++)
261                    {
262                            U.matrixVal[i+U.rows*j]=
263                                    UT.matrixVal[i+UT.rows*j];
264                    }
265            }
266            // solve the minimisation problem using Gaussian
267            // elimination for the upper triangular matrix U.
268            //  This is cheap since it only uses back substitution
269            Vector yu = U / beta;
270
271            // convert back to x
272            Vector kk = Qf * yu;
273            return x0 + kk;
274
275            /* Used for testing gmres convergence
276            Vector residuals1(100);
277            for (int j = 0; j < itsreq; j++ )
278            {
279                    residuals1.matrixVal[j] = residuals.matrixVal[j+1];
280            }
```

```
281          return residuals1;
282          */
283
284 }
```

### C.5.1 Testing GMRES convergence

```
1  #include <iostream>
2  #include  <stdlib.h>
3  # include <cassert>
4  #include "Exception.h"
5  #include "Matrix.h"
6  #include "Vector.h"
7  #include <fstream>
8  int main()
9  {
10         int M = 100;
11         // Create random matrix of size M
12         Matrix A = rand(M)*10.0;
13         // Create RHS TO Ax=b
14         Vector b(M);
15         // Assign 1s to b
16         for (int i = 1; i < M + 1; i++)
17         {
18                 b(i) = 1;
19         }
20         // set gmres parameters
21         int maxits = 1000;;
22         double tol = 1.0e-10;
23         Vector guess = 0.0 * b;
24         // initialise results vector,
25         // gmres has currently been modified to return
26         // the vector of residuals norm(r_k)/norm(b)
27         Vector results;
28         // stream out results to .dat file for MATLAB
29         // processing
30         std::ofstream outfile("gmres_convergence.dat");
31         assert(outfile.is_open());
32         for (int i = 0; i <= 100; i += 10)
33         {
34                 // we consider convergence for A + multiples
35                 // of the identity
36                 results = gmres(A + i * eye(M), b, guess, tol, maxits);
```

40

```
37            outfile << results;
38        }
39        outfile.close();
40
41 }
```

```matlab
1  % MATLAB PROCESSING
2  T = readtable("gmres_convergence.dat");
3  figure
4  hold on;
5  colormap parula
6  for i=1:11
7      ind1 = 100*(i-1)+1;
8      plot(1:100,table2array(T(ind1:ind1+99,:)),"LineWidth",1.5);
9  end
10 X = [.7 .22];
11 Y = [.7 .22];
12 xlabel("Iterations ($k$)",Interpreter="latex",FontSize=30)
13 ylabel("$\frac{||r_k||}{||b||}$  \quad ...
14 ... ",Interpreter="latex",Rotation=360,FontSize=30)
15 annotation('arrow',X,Y);
16 annotation('arrow',X,Y);
17 legend("$0$","$10$","$20$","$30$","$40$","$50$","$60$"...
18 ...,"$70$","$80$","$90$","$100$",Interpreter="latex",fontsize=20)
19 title("Convergence of $A+nI$",Interpreter="latex",FontSize=30)
20 title(legend,"$n$",Interpreter="latex",FontSize=20)
21 txt = '$n$';
22 text(9,.007,txt,Interpreter="latex",FontSize=20)
23 grid on
```

## C.6 Finite Element Solution to Poisson Equation

```cpp
1  #include <iostream>
2  #include <stdlib.h>
3  #include <cassert>
4  #include "Exception.h"
5  #include "Matrix.h"
6  #include "Vector.h"
7  #include <fstream>
8  #include <tuple>
9
10 // 1D basis functions (and derivatives):
```

```cpp
11  double phi1(double x);
12  double phi2(double x);
13  double gradphi1(double x);
14  double gradphi2(double x);
15  // forcing function:
16  double f(double x);
17  double masterToLocal(double x, double xl, double xr);
18  int main()
19  {
20          int meshpts = 1000;
21          Vector mesh = linspace(0, 1, meshpts);
22          Vector dofs = linspace(1, meshpts, meshpts);
23          double h = 1.0 / ((double)meshpts - 1.0);
24
25          // Initialise A matrix (stiffness matrix)
26          Matrix A(meshpts); // square of size meshpts
27          // Initialise F vector (forcing vector)
28          Vector F(meshpts); // size meshpts
29
30          // Initialise DOFS
31          int dof1, dof2;
32
33          // Given basis functions 1-x and x,
34          // the local matrix will be symmetric
35          // so need 3 distinct values:
36          double A11, A12, A22;
37          A11 = (double)
38             1.0/6.0*(1.0/h)*(gradphi1(0.0))*(gradphi1(0.0))
39          + 4.0/6.0*(1.0/h)*(gradphi1(0.5))*(gradphi1(0.5))
40          + 1.0/6.0*(1.0/h)*(gradphi1(1.0))*(gradphi1(1.0));
41
42          A12 = (double)
43             1.0/6.0*(1.0/h)*(gradphi1(0.0))*(gradphi2(0.0))
44          + 4.0/6.0*(1.0/h)*(gradphi1(0.5))*(gradphi2(0.5))
45          + 1.0/6.0*(1.0/h)*(gradphi1(1.0))*(gradphi2(1.0));
46
47          A22 = (double)
48             1.0/6.0*(1.0/h)*(gradphi2(0.0))*(gradphi2(0.0))
49          + 4.0/6.0*(1.0/h)*(gradphi2(0.5))*(gradphi2(0.5))
50          + 1.0/6.0*(1.0/h)*(gradphi2(1.0))*(gradphi2(1.0));
51          // Vector values
52          double F1, F2;
53
```

```
54          // initialise transformed coordinates
55          double x1;
56          double x2;
57          double x3;
58
59          // loop through the elements in the domain,...
60          // ... adding the submatrix to A in the appropriate locations
61          for (int element = 1; element < dofs(meshpts); element++)
62          {
63                  // dofs give the index in A that we add to
64                  dof1 = element;
65                  dof2 = element + 1;
66                  // evaluate Simpsons rule points on the local cells
67                  x1 = masterToLocal(0.0, mesh(dof1), mesh(dof2));
68                  x2 = masterToLocal(0.5, mesh(dof1), mesh(dof2));
69                  x3 = masterToLocal(1.0, mesh(dof1), mesh(dof2));
70                  // evaluate F elements at correct dofs
71                  F1 = (double)1.0 / 6.0 * (h * f(x1) * phi1(0.0)) +
72                          +4.0 / 6.0 * (h * f(x2) * phi1(0.5)) +
73                          +1.0 / 6.0 * (h * f(x3) * phi1(1.0));
74                  F2 = (double)1.0 / 6.0 * (h * f(x1) * phi2(0.0)) +
75                          +4.0 / 6.0 * (h * f(x2) * phi2(0.5)) +
76                          +1.0 / 6.0 * (h * f(x3) * phi2(1.0));
77                  // Matrix and vector assembly
78                  A(dof1, dof1) += A11;
79                  A(dof1, dof2) += A12;
80                  A(dof2, dof1) += A12;
81                  A(dof2, dof2) += A22;
82                  F(dof1) += F1;
83                  F(dof2) += F2;
84          }
85          // Apply boundary conditions
86          A(1, 1) = 1;
87          A(1, 2) = 0;
88          A(meshpts, meshpts) = 1;
89          A(meshpts, meshpts - 1) = 0;
90          F(1) = 0;
91          F(meshpts) = 0;
92
93          // Printing the computed matrix and vector
94          //std::cout << "A: " << A << "\n";
95          //std::cout << "F: " << F << "\n";
96
```

```cpp
 97            // Solving the problem AU = f
 98
 99            // Using Gaussian elimination
100            Vector U_ge(meshpts);
101            U_ge = A / F;
102
103            // Using GMRES
104            Vector U_gm(meshpts);
105            Vector guess = 0 * U_gm;
106            int maxits = meshpts;
107            double tol = 1.0e-8;
108            U_gm = gmres(A, F, guess, tol, maxits);
109
110            // Write results in columns | x | u(x) | to .dat
111            // file for MATLAB processing
112
113
114            // Write Guassian elimination solution for 1000 mesh points
115
116            std::ofstream outfile1("ge_solution_1000.dat");
117            assert(outfile1.is_open());
118            outfile1 << mesh << " " << U_ge << "\n";
119            outfile1.close();
120
121            // Write GMRES solution for 1000 mesh points
122            std::ofstream outfile2("gm_solution_1000.dat");
123            assert(outfile2.is_open());
124            outfile2 << mesh << " " << U_gm << "\n";
125            outfile2.close();
126
127 }
128 // basis functions
129 double phi1(double x){return 1.0 - x;}
130 double phi2(double x){return x;}
131 double gradphi1(double x){return -1.0;}
132 double gradphi2(double x){return 1.0;}
133 // forcing function
134 double f(double x)
135 {
136            return -40.0 / 3.0 * pow(x, 3);
137 }
138 // master to local mapping
139 double masterToLocal(double x, double xl, double xr)
```

```
140  {
141          return xl + x * ( xr − xl );
142  }
```

## C.7  QR algorithm

```
1   \subsubsection{QR decomposition}
2   std::tuple<Matrix, Matrix> qr(const Matrix& A)
3   {
4           // We assume A is square
5           int m = A.rows;
6           // Throw exception if A is non square
7           if (m != A.columns)
8           {
9                   throw Exception(
10                  "Bad_input", "QR_requires_a_sqaure_matrix_input");
11          }
12          // Initialise R as A and Q as mxm identity
13          Matrix R(A);
14          Matrix Q = eye(m);
15
16          // Initialise givens rotataion matrix
17          Matrix givens;
18
19          // Initialise variables for givens rotations
20          double top, bot, div, c, s;
21
22          // the 2x2 givens matrix will be [c s ; -s c ]...
23          // ...where c and s depend on top and bot
24
25          // Loop through each column
26          for (int j = 0; j < m; j++)
27          {
28                  // Loop from the final entry to the first...
29                  // ... subdiagonal entry in each column
30                  for (int i = m - 1; i > j; i--)
31                  {
32                          // define top and bot to be the pair  ...
33                          // ... of adjacent entries on the ...
34                          // ... given column at the given row index
35                          top = R.matrixVal[(i - 1) + m * j];
36                          bot = R.matrixVal[(i)+m * j];
37
38                          // we only want to zero the bottom entry...
39                          // ... if it is non zero
40                          if (bot != 0)
41                          {
```

```
42                                    // clear previous givens matrix,
43                                    // initialise new one
44                                    givens = eye(m);
45                                    // compute givens matrix constants
46                                    div = pow(pow(top,2)+pow(bot,2),0.5);
47                                    c = top / div;
48                                    s = bot / div;
49                                    // assign values to givens matrix
50                                    givens.matrixVal[i+m*i] = c;
51                                    givens.matrixVal[i+m*(i-1)] = -s;
52                                    givens.matrixVal[(i-1)+m*i] = s;
53                                    givens.matrixVal[(i-1)+m*(i-1)] = c;
54                                    // apply givens matrix to zero...
55                                    // ...subdiagonal entries
56                                    R = givens * R;
57                                    // track the orthogonal vectors
58                                    Q = Q * transpose(givens);
59                        }
60                }
61        }
62
63        // FOR TESTING
64        // we can uncomment below to make it more clear when...
65        // ... printing R that it is correct ...
66        // ... since often rounding error adds up.
67        /*
68        double tol = 1e-10;
69        for (int i = 0; i < R.rows * R.columns; i++)
70        {
71                if (fabs(R.matrixVal[i]) < tol)
72                {
73                        R.matrixVal[i] = 0;
74                }
75        }
76        */
77
78        // return tuple (Q,R)
79        return std::make_tuple(Q, R);
80 }
```

### C.7.1   Schur Form

```
1 Matrix eigMatrix(const Matrix& A)
```

```cpp
2    {
3
4            int n = A.rows;
5            // Initialise QR matrices
6            Matrix Q, R;
7            // Initialise Schur form matrix
8            Matrix E(n);
9            // QR already throws for non square input
10           // compute QR
11           std::tie(Q, R) = qr(A);
12           // We impose a hidden maximum iterations for QR
13           int maxits = 1000;
14           // termination condition counter for subdiagonal
15           // zeros
16           int ctr = 0;
17           int i = 0;
18           // tolerence of finding zeros
19           double tol = 1.0e-10;
20           while (i < maxits && ctr < n - 1)
21           {
22                   std::cout << "it:_" << i << "\n";
23                   // mutliply the RxQ
24                   E = R * Q;
25                   //evecs = evecs * Q;
26                   for (int j = 0; j < n-1 ; j++)
27                   {
28                           // check subdiagonal sizes as termination
29                           // condition
30                           if (fabs(E.matrixVal[(j + 1) + n * j]) < tol)
31                           {
32                                   ctr++;
33                           }
34                   }
35                   // compute next QR
36                   std::tie(Q, R) = qr(E);
37                   i++;
38           }
39           // We need to compute one final multiplication
40           E = R * Q;
41           /*
42           * Uncomment for testing - zeros elements close to zero
43           * easier to track
44           for (int i = 0; i < A.rows * A.columns; i++)
```

```
45              {
46                      if  (fabs(E.matrixVal[i]) < tol)
47                      {
48                              E.matrixVal[i] = 0;
49                      }
50              }
51              */
52              return E;
53  }
```

### C.7.2   eig()

```
1  Vector eig(const Matrix& A)
2  {
3          // extracts diagonal elements
4          // from the Schur form of A
5            int n = A.rows;
6            // calculate Schur form
7            Matrix eigenvalueMatrix = eigMatrix(A);
8            Vector eigVec(n);
9            for (int i = 0; i < n; i++)
10           {
11                   // get diagonal elements
12                    eigVec.matrixVal[i]=eigenvalueMatrix.matrixVal[i+n*i];
13           }
14           return eigVec;
15 }
```

## C.8   Eigenvalue problem

```
1  #include <iostream>
2  #include  <stdlib.h>
3  # include <cassert>
4  #include "Exception.h"
5  #include "Matrix.h"
6  #include "Vector.h"
7  #include <fstream>
8  #include <tuple>
9  # define M_PI 3.14159265358979323846
10
11 // 1D basis functions ( and derivatives ):
```

```cpp
12  double phi0(double x);
13  double phi1(double x);
14  double gradphi0(double x);
15  double gradphi1(double x);
16
17  // xi -> x coordinates caluclator
18  double masterToLocal(double x, double xl, double xr);
19
20  // c function
21  double c(double x);
22  int main()
23  {
24          int meshpts = 50;
25          Vector mesh = linspace(0.0, M_PI, meshpts);
26          Vector dofs = linspace(1, meshpts, meshpts);
27          double h = M_PI / ((double)meshpts - 1.0);
28
29          // initialise problem matrices and vector
30          Matrix A(meshpts);
31          Matrix M(meshpts);
32          Vector F(meshpts);
33          // initialise dofs
34          int dof1, dof2;
35
36          // initialise Matrix values
37          double Aii, Aij, Ajj;
38          double Mii, Mij, Mjj;
39
40      // Initialise transformed coordinates
41          double x1, x2, x3;
42
43          // M can be calculated prior to the loop
44          Mii = 1.0 / 6.0 * (h * phi0(0.0) * phi0(0.0))
45                  + 4.0 / 6.0 * (h * phi0(0.5) * phi0(0.5))
46                  + 1.0 / 6.0 * (h * phi0(1.0) * phi0(1.0));
47
48          Mij = 1.0 / 6.0 * (h * phi1(0.0) * phi0(0.0))
49                  + 4.0 / 6.0 * (h * phi1(0.5) * phi0(0.5))
50                  + 1.0 / 6.0 * (h * phi1(1.0) * phi0(1.0));
51
52          Mjj = 1.0 / 6.0 * (h * phi1(0.0) * phi1(0.0))
53                  + 4.0 / 6.0 * (h * phi1(0.5) * phi1(0.5))
54                  + 1.0 / 6.0 * (h * phi1(1.0) * phi1(1.0));
```

50

```
55          for (int element = 1; element < dofs(meshpts); element++)
56          {
57              // assign degrees of freedom - indices for
58              // the global matrix
59                  dof1 = element;
60                  dof2 = element + 1;
61
62              // convert to local coordinates for the calculation
63              // of c(x)
64                  x1 = masterToLocal(0.0, mesh(dof1), mesh(dof2));
65                  x2 = masterToLocal(0.5, mesh(dof1), mesh(dof2));
66                  x3 = masterToLocal(1.0, mesh(dof1), mesh(dof2));
67
68
69              // calculate local elements using Simpsons rule
70                  Aii = (double)
71                  1.0/6.0*((1.0/h)*(gradphi0(0.0))*(gradphi0(0.0))
72                  +h*c(x1)*phi0(0.0)*phi0(0.0))
73                  +4.0/6.0*((1.0/h)*(gradphi0(0.5))*(gradphi0(0.5))
74                  +h*c(x2)*phi0(0.5)*phi0(0.5))
75                  +1.0/6.0*((1.0/h)*(gradphi0(1.0))*(gradphi0(1.0))
76                  +h*c(x3)*phi0(1.0)*phi0(0.0));
77
78                  Aij = (double)
79                  1.0/6.0*((1.0/h)*(gradphi1(0.0))*(gradphi0(0.0))
80                   +h*c(x1)*phi1(0.0)*phi0(0.0))
81                  +4.0/6.0*((1.0/h)*(gradphi1(0.5))*(gradphi0(0.5))
82                  +h*c(x2)*phi1(0.5)*phi0(0.5))
83                  +1.0/6.0*((1.0/h)*(gradphi1(1.0))*(gradphi0(1.0))
84                  +h*c(x3)*phi1(1.0)*phi0(1.0));
85
86                  Ajj = (double)
87                  1.0/6.0*((1.0/h)*(gradphi1(0.0))*(gradphi1(0.0))
88                  +h*c(x1)*phi1(0.0)*phi1(0.0))
89                  +4.0/6.0*((1.0/h)*(gradphi1(0.5))*(gradphi1(0.5))
90                  +h*c(x2)*phi1(0.5)*phi1(0.5))
91                  +1.0/6.0*((1.0/h)*(gradphi1(1.0))*(gradphi1(1.0))
92                  +h*c(x3)*phi1(1.0)*phi1(0.0));
93
94              // append local elements to global matrices
95                  A(dof1, dof1) += Aii;
96                  A(dof1, dof2) += Aij;
97                  A(dof2, dof1) += Aij;
```

```
98                          A( dof2 ,  dof2 )  +=  Ajj ;
99
100                         M( dof1 ,  dof1 )  +=  Mii ;
101                         M( dof1 ,  dof2 )  +=  Mij ;
102                         M( dof2 ,  dof1 )  +=  Mij ;
103                         M( dof2 ,  dof2 )  +=  Mjj ;
104             }
105             // apply boundary conditions
106             A(1 ,  1 )  =  0;
107             A(1 ,  2 )  =  0;
108             A( meshpts ,  meshpts )  =  0;
109             A( meshpts ,  meshpts − 1 )  =  0;
110             M(1 ,  1 )  =  1;
111             M(1 ,  2 )  =  0;
112             M( meshpts ,  meshpts )  =  1;
113             M( meshpts ,  meshpts − 1 )  =  0;
114
115             // compute the matrix we want the eigenvalues of
116             Matrix  B  =  inverse (M)  ∗  A;
117             // print eigenvalues
118             std :: cout  <<  " eigs : ␣"  <<  eig (B)  <<  "\n";
119
120 }
121 // master to local map
122 double masterToLocal(double x, double xl, double xr)
123 {
124             return  xl + x ∗ ( xr − xl );
125 }
126 // c(x) function, we take 0.0 as an example
127 double c(double x)
128 {
129             return  0.0;
130 }
131 // basis functions and their derivatives
132 double phi0(double x)
133 {
134             return  1.0 − x;
135 }
136 double phi1(double x)
137 {
138             return  x;
139 }
140 double gradphi0(double x)
```

```
141  {
142          return −1.0;
143  }
144  double gradphi1(double x)
145  {
146          return 1.0;
147  }
```

## C.9 Miscellaneous Functions

### C.9.1 eye()

```
1  Matrix eye(const int n)
2  {
3          Matrix m(n);
4          for (int i = 0; i < n; i += 1)
5          {
6                  m.matrixVal[i + i * n] = 1;
7          }
8          return m;
9  }
```

### C.9.2 transpose()

```
1   Matrix transpose(const Matrix& m)
2   {
3           Matrix mt(m.columns, m.rows);
4           for (int i = 0; i < mt.rows; i++)
5           {
6                   for (int j = 0; j < mt.columns; j++)
7                   {
8                           mt.matrixVal[i + j * mt.rows] =
9                           m.matrixVal[j + i * m.rows];
10                  }
11          }
12          return mt;
13  }
```

### C.9.3 rand()

```
1  Matrix rand ( int n)
2  {
3          Matrix A(n);
4          srand(time(NULL)); // random seed
5          for (int i = 0; i < n * n; i++)
6          {
7                  A. matrixVal [ i ] =  ((double) rand ()  / (RAND_MAX));
8          }
9          return A;
10 }
```

### C.9.4   inverse()

```
1  Matrix inverse (const Matrix& A)
2  {
3          // uses Gaussian elimination to
4          // construct inverse by solving
5          // Ax=e_i
6          int n = A.rows;
7          Matrix Ainv(n);
8          Vector col;
9          for (int i = 0; i < n; i++)
10         {
11                 Vector b(n);
12                 b. matrixVal [ i ] = 1;
13                 col = A / b;
14                 Ainv.setColumn( i+1, col );
15         }
16         return Ainv;
17 }
```

### C.9.5   convertMatrixToVector()

```
1  Vector convertMatrixToVector (const Matrix& A)
2  {
3          int n = A.rows;
4          int m = A.columns;
5          if (m != 1)
6          {
7                  std :: cout << "BAD_CONVERSION_TO_VECTOR";
8                  return Vector ();
9          }
```

```
10          Vector x(n);
11          for (int i = 0; i < n; i++)
12          {
13                  x.matrixVal[i] = A.matrixVal[i];
14          }
15          return x;
16  }
```

### C.9.6   linspace()

```
1   Vector linspace(double a, double b, int n)
2   {
3           Vector v(n);
4           double length = b − a;
5           double h = length / (n−1);
6           for (int i = 0; i < n ; i++)
7           {
8                   v.matrixVal[i] = a + h * i;
9           }
10          return v;
11  }
```

### C.9.7   norm()

```
1   double norm(Vector v, int p)
2   {
3           double norm_val = 0.0;
4           double temp;
5           for (int i = 0; i < v.rows; i++)
6           {
7                   temp = fabs(v.matrixVal[i]);
8                   norm_val += pow(temp, p);
9           }
10          return pow(norm_val, 1.0 / ((double)(p)));
11  }
```