# C++ for Scientific Computing

Mark Richardson

May 2009

# Contents

# 1   Introduction

The aim of this special topic is to develop C++ code that mirrors the functionality of the scientific computation package, MATLAB. In particular, we will attempt to write code in C++ that enables the user to perform most of the basic, and some of the more complicated matrix and vector operations found in MATLAB. A natural question to ask is *why* should we be interested in replicating the functionality of an existing piece of software? In order to address this, let's briefly look the history of, and what is offered by, respectively, MATLAB and C++.

During the late 1970s, a lecturer at the University of New Mexico, Cleve Moler, wished to give his students access to the linear algebra software libraries, LINPACK and EISPACK, though without them having without them having to learn FORTRAN, the language these software libraries were written in. MATLAB was first developed by Moler in order to allow this access and after going from strength to strength, in 1984 MATLAB was first released as a commercial product. MATLAB has developed a great deal since then and is now widely used in industry and academia for a wide range of problems in scientific computing. It would not be unfair to say that MATLAB is generally considered to be the benchmark tool for computations in linear algebra. MATLAB is popular because of its ease of use. It performs calculations with vectors and matrices, and though MATLAB can perform these computations speedily under certain circumstances, other common programming languages such as C, C++, FORTRAN are generally considered to be faster.

Around the same time at Bell Labs, New Jersey, C++ was being created by Bjarne Stroustrup. C++ was concieved as an extension to the language, 'C', which provided support for classes and consequently, object oriented design. The name of the language, 'C++' is a humorous reference to the command from C for incrementing an integer (to increment an integer, x we would write `x++`). C++ is thus C incremented. C++ has the advantage of having both high and low level features. The object oriented aspect of C++ allows code to be abstracted far away from the machine running the code. However, it is also possible to write code in C++ which accesses the computer's memory directly. C++ has been incredibly successful and is often the language of choice for programmers who wish to develop a piece of professional software.

As powerful as MATLAB undoubtedly is, if we could replicate its functionality in another language, such as C++, we would be free from the shackles of the MATLAB commercial license, and would be able to make any extensions we saw fit. We could even design our own user interface and place it within a separate application. In addition to this, the speed at which code written in C++ executes is incentive enough to make replicating MATLAB commands a worthwhile task.

We will first define a class of matrices. This will be achieved by specifying in our code how the matrix *constructor* should allocate memory to store the details of the matrix, whenever a matrix object is created by the user. In addition to this, we will need to define an algebra for matrix objects, that is we will need to specify exactly what it means (for example) to add, subtract and multiply matrix objects together. Once we have established this basic algebra, we will go on to define a sub-class of vectors. It is then natural to look at trying to implement some of the more sophisticated linear algebra tools. In this project, we will be focussing on solving linear systems of equations. Two methods we shall code for doing this will be the ubiquitous direct method, Gaussian Elimination with partial pivoting, and the Krylov subspace iterative method, Generalised Minimal Residual (GMRES).

## 2  The Matrix Class

Our class of matrices is defined using two files: a header file, *matrix.h*, and C++ code file, *matrix.cpp*. If we wish to run any code that uses this class, we require a *main.cpp* file. Contained within the `main()` scope within this file, we can place the name of other pieces of code that we wish to be executed. The test code for the matrix class is found on a file named, *use_matrix.cpp*, and is referred to within the `main()` scope in the *main.cpp* file as `use_matrix()`.

The header file of a class defines all the *variables* and *methods* that are associated with the class. There are three member variables for our class of matrices; the number of rows of the matrix (`rows`), the number of columns (`columns`), and a pointer to the first entry of a vector of pointers (`**x`) to double precision numbers. This is the standard way of declaring an array of numbers and entries of the matrix can henceforth be accessed using the declaration, `x[i][j]`, where $i$ and $j$ are integers. The methods of the matrix class are functions that act upon variables of a matrix object, or upon the object itself. They are declared in the *matrix.h* header file and the code that implements them is contained in the *matrix.cpp* file.

Each of the variables and methods must be declared as; `public`, `private` or `protected`. The variables of the matrix class are declared as, `protected`, since later we will define a sub-class of vectors which will require to have access to the variables of the matrix class. The methods and constructors are declared as, `public`.

### 2.1  Constructors

Any particular instance of the matrix class is known as a matrix *object*. Whenever such an object is created, a special routine, called the *constructor* is executed. The job of the constructor is to allocate memory for the variables of the object. The matrix class has three constructors; a default constructor, a constructor which takes the matrix dimensions as arguments, and a copy constructor. In addition to this, we also declare a *destructor* which deletes the memory of an object once it goes out of scope. The constructor which creates objects by taking two integer arguments is:

```
matrix::matrix(int no_of_rows, int no_of_columns)
{
    // set the variable values
    rows = no_of_rows;
    columns = no_of_columns;
    xx = new double *[no_of_rows];
    for (int i=0; i<no_of_rows; i++)
    {
        // allocate the memory for the entries of the matrix
        xx[i] = new double[no_of_columns];
    }
    for (int i=0; i<no_of_rows; i++)
    {
        for (int j=0; j<no_of_columns; j++)
        {
            // initialise the entries of the matrix to zero
            xx[i][j]=0.0;
        }
    }
}
```

We see that in lines 4 & 5, integer values are assigned to the *rows* and *columns* fields. Line 6 assigns the memory by creating a vector of pointers (corresponding to the number of rows of the matrix). The loop in lines 7-11 then allocates the number of columns for each row. The final loop initialises the content of the matrix to 0.0. The other constructors and the destructor can be found in the appendix.

## 2.2 Binary Operators

Now that we are able to create matrix objects, we can go about defining what it means to add, subtract, multiply and divide matrices by scalars and other matrices. This is done by *overloading* the binary operators, $+, -, *, /$. Much of this is repetitive, so we shall only describe a few key examples in detail. The full code can be found in the appendix. The following code tells the matrix class what to do if it is asked to add two matrices together.

```
1  matrix operator +(const matrix& A, const matrix& B)
2  {
3      int m = rows(A), n = columns(A), p = rows(B), q = columns(B);
4      if ( m != p || n !=q )
5      {
6          std::cout << "Inconsistent dimensions: Returned first argument";
7          return A;
8      }
9      else
10     {
11         matrix C(m,n);
12         for (int i=0; i<m; i++)
13         {
14             for (int j=0; j<n; j++)
15             {
16                 C.xx[i][j]=A.xx[i][j]+B.xx[i][j];
17             }
18         }
19         return C;
20     }
21 }
```

There is a simple, unsophisticated dimension check in lines 4-8. If the dimensions of the matrix are not identical, then an error message is displayed and the first argument (A) is returned to the user. There are better ways of dealing with such situations (exception handling), but we shall not be discussing them in any great detail here. If the dimensions are consistent, then the code in lines 10-20 is executed. This consists of two nested loops over the rows and columns of the matrix. Quite simply, the sum of the entries of the input matrices (A & B) are added and assigned to the corresponding entry of the output matrix (C). The operator '$-$' is overloaded in a directly analogous way to this. Note that the overloading of '$+$' uses the functions, `rows()` and `columns()`, which are defined as friend functions and return the row and column dimensions. The function `rows()` is defined by:

```
1  int rows(matrix A);           // function prototype
2  int rows(matrix A)
3  {
4      return A.rows;            // Returns the protected field, 'rows'
5  }
```

The unary operators, $+, -$ are defined naturally using similar code to that of the binary

operators. To multiply and divide matrix objects by scalars, we use the following code:

```
1  matrix operator *(const double& p, const matrix& A)
2  {    // Creates a matrix with the same dimensions as A
3      matrix B(A.rows,A.columns);
4          for (int i=0; i<A.rows; i++)
5          {
6              for (int j=0; j<A.columns; j++)
7              {
8                  B.xx[i][j]= p * A.xx[i][j];      // Multiply each entry by p
9              }
10         }
11         return B;
12 }
```

## 2.3  Matrix Multiplication

Multiplying two matrices $\mathbf{A} * \mathbf{B}$, where $A \in \mathbb{R}^{m \times n}$ and $\mathbf{B} \in \mathbb{R}^{n \times p}$, results in a matrix, $\mathbf{C} \in \mathbb{R}^{m \times p}$ where the $(i,j)^{th}$ entry of $\mathbf{C}$ is given by:

$$\mathbf{C}_{i,j} = \sum_{k=1}^{n} \mathbf{A}_{i,k} \mathbf{B}_{k,j}$$

For each $1 \le i \le m$ and $1 \le j \le p$. To overload the '*' operator for two matrix arguments, the following code which implements the above expression is used. It involves three nested loops (one for the rows, one for the columns, and one for the sum for each entry):

```
1  matrix operator *(const matrix& A, const matrix& B)
2  {    // use assertion to check matrix dimensions are consistent
3      assert(A.columns==B.rows);
4
5      // create a result matrix, C with the correct dimensions
6      matrix C(A.rows,B.columns);
7
8      // initialise temp (sum variable)
9      double temp = 0;
10
11         for (int i=0; i < A.rows ; i++)          // for rows (m)
12         {
13             for (int j=0; j < B.columns ; j++)   // for columns (q)
14             {
15                 for (int k=0; k < A.columns ; k++)
16                 {   // dot product step (n sums)
17                     temp = temp + A.xx[i][k]*B.xx[k][j];
18                 }
19
20                 C.xx[i][j]=temp;     // set the C matrix values to temp
21                 temp = 0;            // reset temp
22             }
23         }
24      return C;
25 }
```

Here, we have used a simple assertion to check if the operation is dimensionally consistent. The assertion ends the code immediately if the statement in the brackets is untrue. Note that this code is general enough for us to automatically have defined matrix-vector multiplication, if we define an m-vector as an $m \times 1$ matrix.

5

## 2.4 Achieving MATLAB-like functionality

Now that we have established a basic algebra for matrix objects, we can go about making the code 'feel' more like MATLAB. The two most important ways of doing this are by overloading the assignment operator, '=' and the brackets operator, '()'. In MATLAB if **P** were a matrix, and we were to write `Q=P`, then we would have formed an identical copy of **P**, called **Q**. We need to explicitly tell our matrix class how this should be done. The following code allows us to do this, but only after initially declaring a matrix object that the RHS is being assigned to (this is different to how things work in MATLAB, but unavoidable here).

```
1   matrix& matrix::operator =(const matrix &A)
2   {// Operator returns a matrix equal to the RHS
3
4       // Destruct previous entries
5       for (int i=0; i < rows; i++)
6           {
7               delete[] xx[i];
8           }
9       delete[] xx;
10
11      // Assign new dimensions to be equal to that of the RHS
12      rows = A.rows;
13      columns = A.columns;
14
15      // Allocate the memory as in the constructor
16      xx = new double *[A.rows];
17          for (int i=0; i < A.rows ; i++)
18          {
19              xx[i] = new double[A.columns];
20          }
21
22          // Copy the values across from the RHS
23          for (int i=0; i < A.rows; i++)
24          {
25              for (int j=0; j < A.columns ; j++)
26              {   // Set entries to be the same as the RHS matrix
27                  xx[i][j]=A.xx[i][j];
28              }
29          }
30      return *this;
31  }
```

In all cases, we begin with a matrix object (on the LHS of the '=' sign) that the RHS matrix is being assigned to. Lines 5-9 of the above code delete any data that was there previously (since the user has decided that this object is to be overwritten). This code is very similar to that found in the destructor. Lines 12 and 13 then reset the dimensions of the target matrix to be equal to that of the RHS whilst lines 16-20 reallocate memory just as the constructor did. We now have a matrix with object with the correct dimensions. The last remaining task is to copy across the values from the RHS matrix. This is achieved in lines 23-29.

In MATLAB, it is possible to access and write to elements of a matrix, **A**, by using the command `A(i,j)`, where $i$ and $j$ are the row and column entries of the matrix. To achieve this effect for our matrix class, we use the following code. Note that again we are using a simple check to ensure that the inputted values are within the dimension bounds of the

matrix, though this time we are not using an assertion, merely a warning message.

```
1  double &matrix::operator() (int i, int j)
2  // Allows reference to the entries of a matrix in
3  // the same way as MATLAB. Can call or assign values.
4  {
5      if (i < 1 || j < 1)
6      {
7          std::cout << "Warning: An index may have been too small \n\n";
8      }
9
10     else if (i > rows || j > columns)
11     {
12         std::cout << "Warning: An index may have been too large \n\n";
13     }
14     return xx[i-1][j-1];
15 }
```

In order to debug code effectively, we would like to be able to display the contents of matrices easily. For example, we would like to be able to simply write `std::cout << A;` to display the contents of the matrix, A. This can be achieved by overloading the '`<<`' operator in the following way:

```
1  ostream& operator<< (ostream& output, const matrix &A)
2  {
3      for (int i=0; i<A.rows; i++)
4      {   // work down the rows of the matrix
5
6          for (int j=0; j < A.columns; j++)
7
8          {   // send the entry in each column to output, with a space
9              output << " " << A.xx[i][j];
10         }
11
12         // begin newline in output when the end of the row is reached
13         output << "\n";
14     }
15     output << "\n";
16     return output;
17 }
```

## 2.5  Special Functions and Operations

We will now attempt build up a small library of functions that perform tasks that will be useful to use when we attempt to write code for solving linear systems using either Gaussian Elimination or GMRES. To create an $n \times n$ identity matrix in MATLAB, we would write something like, `I=eye(n)`. We can obtain the same functionality with our matrix class by using the following code:

```
1   matrix eye(int size)
2   {    // Create a temporary matrix with the same dimensions as A
3        matrix I(size,size);
4
5        for (int i=0; i < size; i++)
6        {   // set the values on the diagonal to be 1
7
8               I.xx[i][i] = 1;
9        }
10       return I;
11  }
```

For the partial pivoting steps in the Gaussian Elimination algorithm (which will be explained in greater detail later), we will need to identify the pivot for a given column of a matrix (this is the largest entry that is either on the diagonal or below the diagonal for a given columns). The code that achieves this takes a matrix and an integer (the column number you wish to find the pivot for) as inputs and returns the row number corresponding to the largest diagonal/sub-diagonal value for that column:

```
1   int find_pivot(matrix A, int column)
2   {
3       // Initialise maxval to be diagonal entry in column, 'column'
4       double maxval=fabs(A(column,column));
5       // Initialise rowval to be column
6       int rowval=column;
7
8           for (int i=column+1; i ≤ A.rows; i++)
9           {
10              if ( fabs(A(i,column)) > maxval)
11              {   // Update maxval and rowval if bigger than previous maxval
12                  maxval = fabs(A(i,column));
13                  rowval = i;
14              }
15          }
16      return rowval;
17  }
```

The code works by initialising the variable 'maxval' to be the diagonal entry of the column that is specified by the user as one of the input arguments. The quantity 'rowval' is initialised to be the same as the input argument, 'column'. In the 'for' loop in lines 8-15, the code cycles down the entries below the sub-diagonal and asks if each entry is larger than 'maxval'. If it is, then 'maxval' and 'rowval' are updated. The returned integer value is thus the row corresponding to the largest diagonal/sub-diagonal entry.

Once this information has been obtained, it needs to be used to obtain a permutation matrix which will switch the two rows in question. A function which accepts two integers (representing the two rows to be swapped) and which returns the appropriate permutation matrix is defined using the code below.

We simply create an appropriately sized identity matrix and manually swap over rows $i$ and $j$ by setting the diagonal entries on those rows to zero and the $(i, j)$ and $(j, i)$ entries to one:

```
1  matrix permute_r(int n, int i, int j)
2  {
3      // Create nxn identity matrix
4      matrix I=eye(n);
5
6      // set to zero the diagonal entries in the given rows
7      I(i,i)=0;
8      I(j,j)=0;
9
10         // set the appropriate values to be 1
11      I(i,j)=1;
12      I(j,i)=1;
13
14      return I;
15  }
```

A function that will be required by the GMRES routine is a function that takes a matrix and two integers as inputs and returns a re-sized matrix with dimensions corresponding to the integer inputs. Any entries that are shared by the inputted and outputted matrix are copied across, whilst any entries that are not shared are either lost (if the new dimensions are smaller than the previous dimensions) or set to zero (if the converse is true). This resizing operation can be achieved with the following code:

```
1  matrix resize(matrix A, int m, int n)
2  {
3      int p,q;
4      matrix Mout(m,n);
5
6      if (m≤A.rows)        // set p as the lowest of the two row dimensions
7      {
8          p=m;
9      }
10     else
11     {
12         p=A.rows;
13     }
14
15     if (n≤A.columns)    // set q as the lowest of the two column dimensions
16     {
17         q=n;
18     }
19     else
20     {
21         q=A.columns;
22     }
23
24     for (int i=1; i≤p; i++)      // loop across the smallest row dimension
25     {
26         for (int j=1; j≤q; j++)  // loop across the smallest column dimension
27         {
28             Mout(i,j) = A(i,j);  // copy across the appropriate values
29         }
30     }
31
32     return Mout;
33  }
```

Lastly, we define the method for computing the transpose of a matrix. This will be used in the GMRES code for computing inner products. In MATLAB, once a matrix object, **X** exists, its transpose is computed by the command, X'. It is not as straightforward to allow this command mean transpose in C++, so we will assign the ' ' symbol placed before the object to represent transpose. It is therefore reasonable to class this as a unary operator. We simply create a matrix object with the dimensions reverse and loop through setting values using the transpose relation, $\mathbf{A}_{i,j}^T = \mathbf{A}_{j,i}$.

```
1  matrix operator ¬(const matrix& A)
2  {
3          // Create a temporary matrix with reversed dimensions
4      matrix B(A.columns,A.rows);
5
6          // Set the entries of B to be the same as those in A
7      for (int i=0; i < A.columns; i++)
8      {
9          for (int j=0; j < A.rows; j++)
10             {
11                 B.xx[i][j] = A.xx[j][i];
12             }
13     }
14     return B;
15 }
```

# 3   The Vector Sub-Class

We can consider the set of all vectors to be a subset of the set of all matrices. For consistency with the rules of multiplication defined in the matrix class, we must define vectors to have column dimension 1. This convention of course does not conflict with any notions that we naturally have about the dimensions of vectors. All our previous algebra that was defined for matrices can be applied to vector objects. The vector class is instructed to inherit from the matrix class using the following code in the header file:

```cpp
// 'vector' inherits from 'matrix' in a public way
class vector: public matrix
{
        ...
}
```

As in MATLAB, we would like to construct a vector object with the call, `vector v(n)`, to create a column vector of length $n$. To this we define the vector constructor to run the matrix constructor with $n$ as the first argument and 1 as the second argument.

```cpp
vector::vector(int no_of_elements)
// run the matrix constructor
: matrix(no_of_elements,1)
{ }
```

The default vector constructor simply runs the default matrix constructor. Any instance of the vector class can be passed as an argument into a matrix method and we hence retain all the functionality of the matrix class with vector objects. Mathematically, there is no real *need* to have defined this subclass, the reason it has been added is to make the interface more natural to the user. It also allows us to overload the brackets operator for vectors so that elements can be accessed using the call, `v(i)`, where $i$ is some entry of the vector, `v`. The code for doing this is almost identical to the corresponding matrix method, so we shall not include it here, though it can, along with all the other code from this project, be found in the appendix.

A problem that arises when defining the vector sub-class in this way is that the methods written for the matrix class always return matrix objects. In other words, though a user may pass a vector object as an argument into a matrix method, if that method returns an object, then that object will be a matrix, not a vector. This did not seem to be a problem at first, however when writing the code for GMRES, which uses inner products and other vector operations, it became frustrating to program. The solution to this was to replicate (or overload) the necessary methods for the vector subclass. Thus, the binary, unary and assignment operators have been overloaded and for simple vector operations, the returned object is also a vector. Again, the code for these methods is very similar to the corresponding matrix versions, so they will not be included here.

There are only a few extra methods which are unique to the matrix class, `mat2vec()` and `norm()`. They have both been added to be used in the GMRES code. The first method, `mat2vec()`, takes a matrix object as an argument and returns a vector object. The reason for the existence of such a method is similar to the reasons described above for overloading the binary, unary and assignment operators for the vector class. It is needed if, for example, we wish the result of a matrix-vector multiplication to be a vector. The code accepts a matrix

11

input (which will usually be an $m \times 1$ matrix), creates a vector object with the appropriate dimensions, and copies across the entries. A matrix with more than one column can be passed into the method. In this case, the returned object will just be the first columns of that matrix. The code for this method is:

```
1  vector mat2vec(matrix A)
2  {
3      // create vector with same no of rows as A and one column
4      vector v(rows(A));
5
6      for (int i=1; i ≤ rows(A); i++)
7      {   // copy only the values in the first column
8          v(i)=A(i,1);
9      }
10      return v;
11  }
```

The last function that we will need in order to program the linear solvers is a way to compute the p-norm of a vector. In the header file, the function is declared in the following way:

```
1  friend double norm(vector v, int p=2);
```

The function takes two arguments as inputs; the vector we wish to find the p-norm of, and an integer, $p$. The statement of second argument, `int p=2`, allows the function to be called with just the first argument, in which case the default value, $p = 2$, is used. The following code is adapted from the similar method given in the C++ for scientific computing course on the MSc Mathematical Modelling and Scientific Computing, 2009.

```
1  double norm(vector v, int p)
2  {
3      // define variables and initialise sum
4      double temp, value, sum = 0.0;
5
6      for (int i=1; i ≤ rows(v); i++)
7      {   // floating point absolute value
8          temp = fabs(v(i));
9          sum += pow(temp, p);
10      }
11
12      value = pow(sum, 1.0/((double)(p)));
13
14      return value;
15  }
```

# 4 Gaussian Elimination with Partial Pivoting

The ubiquitous direct method for solving the linear system, $\mathbf{Ax} = \mathbf{b}$, is *Gaussian Elimination*. This method is so important that it is usually first learned in a certain form by secondary school children when they try to solve (usually two or three variable) simultaneous equations.

For the general $n \times n$ case, the goal is to perform a series of operations on the linear system so that we end up with the equation, $\mathbf{Ux} = \mathbf{y}$, where $\mathbf{U}$ is an upper triangular matrix and $\mathbf{y}$ is a known vector. This system can then be solved easily using back-substitution to obtain $\mathbf{x}$. This process can be interpreted in terms of a series of matrix multiplications applied to the left of the original equation, $\mathbf{Ax} = \mathbf{b}$. The basic Gaussian Elimination algorithm begins with the first column of the matrix and attempts to set to zero all the sub-diagonal entries of that column by multiplying them by a suitable multiple of the number in the diagonal entry. Once this has been achieved for that column, the same is applied to the next column, and so on until the matrix has been reduced to an upper-triangular form. Formally stated, the basic GE algorithm is:

> **for** j = 1,2,...,n-1 **do**
>    **for** i = j+1,j+2,...,n **do**
>       calculate multiplier $l_{ij} = \frac{a_{ij}}{a_{jj}}$
>       row i $\leftarrow$ row i - $l_{ij} \times$ row j
>    **end for**
> **end for**

Each outer loop step (over $j$) can be represented as a matrix multiplication. To demonstrate this, consider the simple $4 \times 4$ linear system, $\mathbf{Ax} = \mathbf{b}$:

$$
\begin{pmatrix}
8 & 3 & 4 & 6 \\
2 & 4 & -1 & 2 \\
1 & -4 & -6 & -9 \\
4 & -3 & -4 & -5
\end{pmatrix}
\begin{pmatrix}
x_1 \\ x_2 \\ x_3 \\ x_4
\end{pmatrix}
=
\begin{pmatrix}
1 \\ 2 \\ 3 \\ 4
\end{pmatrix}
$$

We can perform the first step of the GE algorithm by forming an identity matrix and placing the row multiplier multiplied by $-1$ on the appropriate sub-diagonal entry. It is easy to see what the multipliers are in this case. Multiplying $\mathbf{A}$ on the left by this matrix produces:

$$
\begin{pmatrix}
1 & 0 & 0 & 0 \\
-\frac{1}{4} & 1 & 0 & 0 \\
-\frac{1}{8} & 0 & 1 & 0 \\
-\frac{1}{2} & 0 & 0 & 1
\end{pmatrix}
\begin{pmatrix}
8 & 3 & 4 & 6 \\
2 & 4 & -1 & 2 \\
1 & -4 & -6 & -9 \\
4 & -3 & -4 & -5
\end{pmatrix}
=
\begin{pmatrix}
8 & 3 & 4 & 6 \\
0 & 3.25 & -2 & 0.5 \\
0 & -4.375 & -6.5 & -9.75 \\
0 & -4.5 & -6 & -8
\end{pmatrix}
$$

This concept can be implemented successively for each column of the matrix. In order to make the sub-diagonal entries zero, we multiply each matrix on the left by the new matrix for each column. This process will eventually produce an upper triangular matrix, and indeed we will actually obtain an LU factorisation of A. This step in the GE algorithm is implemented with the code from lines 13-21, in the implementation given below. For each column of the matrix, we first create a new identity matrix (line 13). The `for` loop in lines 15-20 then calculates and stores the multiplier in the appropriate sub-diagonal entry. Line 21 then computes the matrix multiplication for that step of the loop. The implementation for GE w. partial pivoting is given overleaf:

```
1   matrix operator /(const matrix& b, const matrix& A)
2   {
3       int n = A.rows;
4       // create empty matrices, P & L, U & Atemp
5       matrix P, L, Utemp = eye(n), Atemp=A;
6
7       for (int j=1; j < n ; j++)
8       {   // create permutation matrix, P
9           P = permute_r(n,find_pivot(Atemp,j),j);
10          // update U & Atemp
11          Utemp = P*Utemp;
12          Atemp = Utemp*A;
13          L = eye(n);
14
15          for (int i=j+1; i <= n ; i++)
16          {   // check for division by zero
17              assert(fabs(Atemp(j,j))>1.0e-015);
18              // compute multiplier and store in sub-diagonal entry of L
19              L(i,j)= -Atemp(i,j)/Atemp(j,j);
20          }
21          Utemp = L*Utemp;
22          Atemp = Utemp*A;
23      }
24
25      matrix U = Utemp*A;
26      matrix y = Utemp*b;
27      matrix x(n,1);                          // Create result vector
28
29      // Solve Ux=y by back substitution:
30      // Compute last entry of vector x (first step in back subs)
31      x(n,1)=y(n,1)/U(n,n);
32      double temp = 0;
33
34      for (int i=n-1; i>=1; i--)
35      {
36          temp = y(i,1);
37          for (int j=n; j>i; j--)
38          {
39              temp = temp - U(i,j)*x(j,1);
40          }
41          x(i,1)=temp/U(i,i);
42      }
43      return x;
44  }
```

Working with exact arithmetic, the basic GE algorithm (assuming that we have 'enough' time and computer power available) will be able to successfully solve all nonsingular linear systems. However, we live in a world of floating point computations. It is very well known that basic GE is in general unstable when implemented on a computer. This is due to the existence floating point rounding errors which can often propagate through the algorithm.

The fix for this issue is to use 'partial pivoting'. This strategy involves permuting the rows of the matrix at each stage of the computation. We ensure that for each column that is currently being manipulated, the largest diagonal or sub-diagonal entry (in absolute terms) is moved to the diagonal entry. When the largest number is already on the subdiagonal, no change is made. As was the case before, this goal can be accomplished using matrix products, in particular, we can use permutation matrices. To illustrate, let's perform this step of the algorithm for the above $4 \times 4$ system.

Recall that we had already achieved our goal for the first column. Our matrix is now:

$$\begin{pmatrix} 8 & 3 & 4 & 6 \\ 0 & 3.25 & -2 & 0.5 \\ 0 & -4.375 & -6.5 & -9.75 \\ 0 & -4.5 & -6 & -8 \end{pmatrix}$$

Our attention is now focussed on the diagonal and sub-diagonal entries of the second column of this matrix. We see that the largest value in absolute terms (-4.5) is in the $(4, 2)$ entry of the matrix. The goal is to move this number to the diagonal of the same column $(2, 2)$. We accomplish this using a permutation matrix to interchange the second and fourth rows:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} 8 & 3 & 4 & 6 \\ 0 & 3.25 & -2 & 0.5 \\ 0 & -4.375 & -6.5 & -9.75 \\ 0 & -4.5 & -6 & -8 \end{pmatrix} = \begin{pmatrix} 8 & 3 & 4 & 6 \\ 0 & -4.5 & -6 & -8 \\ 0 & -4.375 & -6.5 & -9.75 \\ 0 & 3.25 & -2 & 0.5 \end{pmatrix}$$

This part of the algorithm in implemented in lines 8-11 of the code above. We utilise the two methods, `permute r(.,.,.)` and `find pivot(.,.)` described earlier to find the pivot in the given column and then create the corresponding permutation matrix. The actual matrix multiplication occurs in line 11.

Throughout this code, the temporary matrix objects `Utemp` and `Atemp`. `Utemp` is the matrix result of all the matrix operations that have been carried out up to a given point in the algorithm. `Atemp` is `Utemp` multiplied by A for a given point in the algorithm (see line 22). The reason for needing `Atemp` is to allow us to calculate the (step specific) multipliers for the elimination step (line 19).

Once this transformation to an upper-triangular matrix has taken place, the last task the algorithm accomplishes is to solve the triangular system using back substitution. This is found in lines 31-42 of the above code. We start at the bottom of the matrix and solve by working upwards, substituting in all our previously computed unknowns to work out the next unknown.

# 5 GMRES

The work required to implement Gaussian Elimination (even without partial pivoting) is around $\frac{2}{3}n^3$ flops[1]. For large systems, we may look to an iterative solver as an alternative as for certain classes of matrix (e.g. sparse), as they can often converge to an acceptable level of precision with far less work.

During this section, we shall implement the now famous Generalised Minimal Residual algorithm (GMRES). The paper describing this method was first published by Youcef Saad and Martin Schulz in 1986[2]. It has since been one of the most cited papers in all of applied mathematics.

GMRES is a Krylov subspace method which seeks to minimise the residual ($\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{X}$) when measured using the vector 2-norm. Given an initial guess vector, $\mathbf{x_0}$, and corresponding initial residual, $\mathbf{r_0}$, the $k^{th}$ Krylov subspace, $K_k(\mathbf{A}, \mathbf{r_0})$ is defined as:

$$K_k(\mathbf{A}, \mathbf{r_0}) = \text{span}\{\mathbf{r_0}, \mathbf{A}\mathbf{r_0}, \mathbf{A^2}\mathbf{r_0}, \ldots, \mathbf{A^{k-1}}\mathbf{r_0}\}$$

For the $k^{th}$ iterate, Krylov subspace methods attempt to find:

$$\mathbf{x}_k \in \mathbf{x}_0 + K_k(\mathbf{A}, \mathbf{r_0})$$

In general, $\{\mathbf{r_0}, \mathbf{A}\mathbf{r_0}, \mathbf{A^2}\mathbf{r_0}, \ldots, \mathbf{A^{k-1}}\mathbf{r_0}\}$ is particularly bad basis for the Krylov subspace, $K_k(\mathbf{A}, \mathbf{r_0})$, so usually some kind of orthogonalisation process is required. In the case of GMRES, the *Arnoldi* iteration is used (for MINRES, Lanczos). The Arnoldi iteration generates an orthonormal basis for the $k^{th}$ Krylov subspace, using a Gramm-Schmidt style iteration. The Arnoldi iteration is outlined in the following pseudo-code:

> Guess $\mathbf{x_0}$, $\quad\quad$ $\mathbf{r_0} = \mathbf{b} - \mathbf{A}\mathbf{x_0}$, $\quad\quad$ set $\mathbf{v_1} = \frac{\mathbf{r_0}}{||\mathbf{r_0}||_2}$
> **for** l = 1,2,... **do**
> $\quad\quad$ $\mathbf{w} = \mathbf{A}\mathbf{v_l}$
> $\quad$ **for** j = 1,2,...,l **do**
> $\quad\quad\quad$ $h_{j,l} = \mathbf{v_j^T}\mathbf{w}$
> $\quad\quad\quad$ $\mathbf{w} = \mathbf{w} - h_{j,l}\mathbf{v_j}$
> $\quad$ **end for**
> $\quad\quad$ $h_{l+1,l} = ||\mathbf{w}||_2$
> $\quad\quad$ $\mathbf{v_{l+1}} = \frac{\mathbf{w}}{h_{l+1,l}}$
> **end for**

This process can also be represented in matrix form. Defining the $n \times k$ matrix $\mathbf{V_k} = [\mathbf{v_1}, \mathbf{v_2}, ..., \mathbf{v_k}]$, and the $(k+1) \times k$ upper Hessenberg matrix, $\tilde{\mathbf{H}}$, who's entries are created in the Arnoldi process as:

$$\tilde{\mathbf{H}_\mathbf{k}} = \begin{pmatrix} h_{1,1} & h_{1,2} & h_{1,3} & \cdots & h_{1,k-1} & h_{1,k} \\ h_{2,1} & h_{2,2} & h_{2,3} & \cdots & h_{3,k-1} & h_{2,k} \\ 0 & h_{3,2} & h_{3,3} & \cdots & h_{3,k-1} & h_{3,k} \\ \vdots & \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & h_{k,k-1} & h_{k,k} \\ 0 & 0 & 0 & \cdots & 0 & h_{k+1,k} \end{pmatrix}$$

We can then write the Arnoldi process as:

$$\mathbf{A}\mathbf{V_k} = \mathbf{V_{k+1}}\tilde{\mathbf{H}_\mathbf{k}}$$

16

Note that $\mathbf{V_{k+1}}$ is simply the matrix of orthogonal basis vectors for the $(k+1)^{th}$ Krylov subspace. Returning to our task, find $\mathbf{x}_k \in \mathbf{x_0} + K_k(\mathbf{A}, \mathbf{r_0})$, we note that since the columns of $\mathbf{V_k}$ are, by construction, an orthonormal basis for $K_k(\mathbf{A}, \mathbf{r_0})$, this is equivalent to writing:

$$\mathbf{x_k} = \mathbf{x_0} + \mathbf{V_k y} \qquad \text{for some coefficient vector,} \qquad \mathbf{y} \in \mathbb{R}^k$$

$$\begin{aligned}
\text{Now, since} \qquad \mathbf{x} - \mathbf{x_k} &= \mathbf{x} - \mathbf{x_0} + \mathbf{V_k y} \\
\Rightarrow \qquad \mathbf{A}(\mathbf{x} - \mathbf{x_k}) &= \mathbf{A}(\mathbf{x} - \mathbf{x_0} - \mathbf{V_k y}) \\
\Rightarrow \qquad \mathbf{b} - \mathbf{A}\mathbf{x_k} &= \mathbf{b} - \mathbf{A}\mathbf{x_0} - \mathbf{A}\mathbf{V_k y} \\
\Rightarrow \qquad \mathbf{r_k} &= \mathbf{r_0} - \mathbf{A}\mathbf{V_k y}
\end{aligned}$$

Now, let us suppose that we wish to minimise the $k^{th}$ residual in the 2-norm. Our problem is thus: find $\mathbf{y} \in \mathbb{R}^k$ such that $||\mathbf{r_0} - \mathbf{A}\mathbf{V_k y}||_2$ is minimised. From the first step in the Arnoldi process, we defined the first orthonormal basis vector, $\mathbf{v_1}$ using $\mathbf{r_0} = ||\mathbf{r_0}||_2\mathbf{v_1}$. We now make the innocuous seeming observation that:

$$\mathbf{v_1} = \mathbf{V_{k+1}}\mathbf{e_1} \qquad \text{where} \qquad \mathbf{e_1} = (1,0,0,...,0)^T \in \mathbb{R}^{k+1}$$

$$\text{therefore} \qquad \mathbf{r_0} = ||\mathbf{r_0}||\, \mathbf{V_{k+1}}\mathbf{e_1}$$

Recalling also the relation, $\mathbf{A}\mathbf{V_k} = \mathbf{V_{k+1}}\tilde{\mathbf{H}}_\mathbf{k}$, we have:

$$\begin{aligned}
||\mathbf{r_k}||_2 = ||\mathbf{r_0} - \mathbf{A}\mathbf{V_k y}||_2 &= \left|\left| ||\mathbf{r_0}||\,\mathbf{V_{k+1}}\mathbf{e_1} - \mathbf{V_{k+1}}\tilde{\mathbf{H}}_\mathbf{k}\mathbf{y} \right|\right|_2 \\
&= \left|\left| \mathbf{V_{k+1}}\left( ||\mathbf{r_0}||_2\,\mathbf{e_1} - \tilde{\mathbf{H}}_\mathbf{k}\mathbf{y} \right) \right|\right|_2 \\
&= \left|\left| ||\mathbf{r_0}||_2\,\mathbf{e_1} - \tilde{\mathbf{H}}_\mathbf{k}\mathbf{y} \right|\right|_2
\end{aligned}$$

The last line is true because the orthogonality of the columns of $\mathbf{V_{k+1}}$ means that mutliplying by that matrix does not affect the vector 2-norm. Let us remind ourselves of the dimensions of each of the quantities involved. $||\mathbf{r_0}||_2\,\mathbf{e_1}$ is a vector of length $k+1$, $\tilde{\mathbf{H}}_\mathbf{k}$ is a $(k+1) \times k$ matrix, and $\mathbf{y}$ is a vector of length $k$. This is therefore a (Hessenberg) linear least squares problem. The standard way of solving such problems is by using Givens rotation matrices, (since Givens rotation matrices are orthogonal, this in essence leads to a QR factorisation of the Hessenberg matrix).

We are now ready to state the full GMRES algorithm:

$$\text{For arbitrary } \mathbf{x_0}, \qquad \mathbf{r_0} = \mathbf{b} - \mathbf{A}\mathbf{x_0}, \qquad \text{set } \mathbf{v_1} = \frac{\mathbf{r_0}}{||\mathbf{r_0}||_2}$$

**for** k = 1,2,... **do**
    do step k of the Arnoldi process
    (this gives us a new vector, $\mathbf{v_{k+1}}$ and a new last column of the matrix, $\tilde{\mathbf{H}}_\mathbf{k}$
    solve the Hessenberg linear least squares problem:
        find $\mathbf{y} \in \mathbb{R}^n$ s.t. $\left|\left| ||\mathbf{r_0}||_2\,\mathbf{e_1} - \tilde{\mathbf{H}}_\mathbf{k}\mathbf{y} \right|\right|_2$ is minimised
    then, $\mathbf{x_k} = \mathbf{x_0} + \mathbf{V_k y}$
**end for**

The C++ code to implement GMRES is given below. An explanation of the lines of code can be found after this.

```
1  vector GMRES(matrix A, matrix b, matrix x0, double tol)
2  {
3      // determine initial residual, r0 in vector form
4      vector r0 = mat2vec(b − A*x0);
5
6      // need this in least square part later
7      double normr0 = norm(r0);
8      double residual=1.0;
9      vector v= r0/normr0;
10
11      // declare and intitialise variables
12      int k=1;
13      matrix J, Jtotal=eye(2),H(1,1), Htemp, HH;
14      matrix bb(1,1), c, cc, tempMat, V, Vold, hNewCol;
15      vector w, vj(rows(v));
16
17      bb(1,1)=normr0;
18
19      // initialise matrix V (matrix of orthogonal basis vectors)
20      V=v;
21
22      while (residual>tol)
23      {
24          H=resize(H,k+1,k);
25
26          // Arnoldi steps (using Gram−Schmidt process)
27          w = mat2vec(A*v);
28
29              for (int j=1; j≤k; j++)
30              {
31                  for (int i=1; i≤rows(V); i++)
32                  {
33                      // set the vector vj to be jth column of V
34                      vj(i)=V(i,j);
35                  }
36
37                  // the next two lines calculate the inner product
38                  tempMat = ¬vj*w;
39                  H(j,k)= tempMat(1,1);
40                  w = w − H(j,k)*vj;
41              }
42
43          H(k+1,k)=norm(w);
44
45          v=w/H(k+1,k);
46
47          // append an additional column to matrix V
48          V=resize(V,rows(V),k+1);
49
50          for (int i=1; i≤rows(V); i++)
51          {
52              // copy entries of v to new column of V
53              V(i,k+1)=v(i);
54          }
55
56        ///////////////// Least squares step /////////////////////
57
58          if (k==1)
59          {   // First pass through, Htemp=H
60              Htemp=H;
61          }
```

```
62
63          else
64          {
65              // for subsequent passes, Htemp=Jtotal*H
66              Jtotal=resize(Jtotal,k+1,k+1);
67              Jtotal(k+1,k+1)=1;
68              Htemp=Jtotal*H;
69          }
70
71      // form next Givens rotation matrix
72      J = eye(k-1);
73      J = resize(J,k+1,k+1);
74
75      // set values to eliminate the h(k+1,k) entry
76      J(k,k)=Htemp(k,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
77      J(k,k+1)=Htemp(k+1,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
78      J(k+1,k)=-Htemp(k+1,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
79      J(k+1,k+1)=Htemp(k,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
80
81          // combine together with previous Givens rotations
82          Jtotal=J*Jtotal;
83          HH=Jtotal*H;
84
85          bb=resize(bb,k+1,1);
86          c=Jtotal*bb;
87
88          residual=fabs(c(k+1,1));
89
90          k++;
91      }
92
93      std::cout<< "GMRES iteration converged in " << k-1 << " steps\n\n";
94
95      // Extract upper triangular square matrix
96      HH=resize(HH,rows(HH)-1,columns(HH));
97      cc=resize(c,rows(HH),1);
98
99      // solve linear system
100     matrix yy = cc/HH;
101
102     vector y = mat2vec(yy);
103
104     // chop the newest column off of matrix V
105     V=resize(V,rows(V),columns(V)-1);
106
107     vector x= mat2vec(x0+V*y);
108     return x;
109 }
```

The code takes in three primary arguments: the matrix, $\mathbf{A}$, the RHS vector , $\mathbf{b}$, and the initial guess vector, $\mathbf{x_0}$. In addition, there is an optional argument, `tol`, or tolerance, which is set to a default value of $1e-6$, as in MATLAB.

We begin, in line 4 by computing the initial residual vector, $\mathbf{r0}$. Note that we employ `mat2vec()` to ensure that the result is a vector. Line 7 declares a double number `normr0` which is equal to $||\mathbf{r_0}||_2$. We also initialise the double number, `residual=1.0`, in order to enter the subsequent while loop. Line 9 declares the first orthogonal basis vector, $\mathbf{v_1}$. We will store the vectors, $\mathbf{v_i}$, in the matrix, $\mathbf{V}$, so that as we pass through the algorithm, we can just re-use the same vector, `v`, for the current $\mathbf{v_i}$.

Lines 12-15 declare and initialise all the other matrices and vectors that the algorithm uses. Line 17 sets the first entry of the vector, `bb` to be equal to `normr0`, representing the term $||\mathbf{r_0}||_2\mathbf{e_1}$. Note that we can increment the size of this vector as required during the course of the algorithm by using the command, `resize()`. Line 20 initialises the matrix, $\mathbf{V}$, storing $\mathbf{v_1}$ as its first column.

We now enter the main GMRES loop which utilises: `while (residual>tol)`. Line 24 uses the command, `resize()`, to set the upper Hessenberg matrix, $\mathbf{H_k}$, to the correct dimensions for the current pass (recall that `resize()` retains the previous entries if it increases the dimensions). Lines 27-45 simply implement the Arnoldi algorithm as described previously and hence computes the entries for the newest ($k_{th}$) column of $\mathbf{H_k}$. A small problem was that when calculating the inner product in lines 38 & 39, the computation is a matrix-matrix multiplication. Therefore, the returned object is a matrix, not a double. We have to manually extract the $(1, 1)$ entry of the resultant matrix, `tempmat`, in order to get the required result. This could well have been placed into a function, but since there was only two lines of code involved, I decided that doing it this way would be easier. Line 45 gives the latest basis vector, $\mathbf{v_{k+1}}$ generated by the Arnoldi process. Lines 48-54 then add an additional column onto the matrix, $\mathbf{V}$ and stores the new vector in it by copying across the values.

We now proceed to the least squares step. Recall that this involves converting the Hessenberg matrix to an upper triangular matrix and then solving the truncated linear system (if the residual is small enough). In order to make the Hessenberg matrix upper triangular, we use Givens rotation matrices. During this part of the algorithm, we will use three additional matrices; J, the Givens rotation matrix for the *current* step; Jtotal, the running product of all the previous Givens rotation matrices; and `Htemp=Jtotal*H`, which is the effect that all the previous rotations have upon the current Hessenberg matrix. Proceeding in this way allows us to only require one Givens rotation matrix at each step. The computation of `Htemp` allows us to work out what the next Givens rotation matrix should be (`Htemp` is always 'almost' upper-triangular - the only nonzero sub-diagonal entry is $\mathbf{Htemp}_{k+1,k}$. Lines 58-79 accomplish this task.

Once we have found the latest Givens rotation matrix, we compute the `Jtotal` in line 82 and then form the upper triangular matrix, `HH=Jtotal*H`, in line 83. This is effectively our QR factorisation of the matrix, $\mathbf{H_k}$. To elucidate: $\mathbf{Jtotal}$ is orthogonal (the product of $k$ Givens rotation matrices) and $\mathbf{HH}$ is upper triangular. We have therefore factorised $\mathbf{H}$ in the following way:

$$\mathbf{H} = (\mathbf{Jtotal^T})(\mathbf{HH})$$

$$\text{i.e.} \quad \mathbf{H} = \mathbf{QR} \quad \text{where} \quad \mathbf{Q} = \mathbf{Jtotal^T} \quad \text{and} \quad \mathbf{R} = \mathbf{HH}$$

Returning to the formal statement of the least squares problem, we have, for step $k$:

$$\begin{aligned}
\left|\left|\,||\mathbf{r_0}||_2\,\mathbf{e_1} - \mathbf{\tilde{H}_k y}\right|\right|_2 &= \left|\left|\,||\mathbf{r_0}||_2\,\mathbf{e_1} - \mathbf{QRy}\right|\right|_2 \\
&= \left|\left|\mathbf{QQ^T}\,||\mathbf{r_0}||_2\,\mathbf{e_1} - \mathbf{QRy}\right|\right|_2 \\
&= \left|\left|\mathbf{Q}\left(\mathbf{Ry} - \mathbf{Q^T}\,||\mathbf{r_0}||_2\,\mathbf{e_1}\right)\right|\right|_2 \\
&= \left|\left|\mathbf{Ry} - \mathbf{Q^T}\,||\mathbf{r_0}||_2\,\mathbf{e_1}\right|\right|_2
\end{aligned}$$

Framing this last statement in terms of the names used in the C++ code, we have:

$$\left|\left|\mathbf{Ry} - \mathbf{Q^T}\,||\mathbf{r_0}||_2\,\mathbf{e_1}\right|\right|_2 = ||(\mathbf{HH}) * (\mathbf{y}) - (\mathbf{Jtotal}) * (\mathbf{bb})||_2$$

Line 86 in the code makes the definition `c=Jtotal*bb`. Since we have a least squares problem with, effectively a $k \times k$ matrix, to find a $k$ dimensional solution, $\mathbf{y}$ with a RHS vector, $\mathbf{c}$ of dimension $k+1$, we know that the $k^{th}$ residual will be the $k+1^{th}$ entry of the vector, $\mathbf{c}$. This is what is computed in line 88. Note that we do not at this stage solve the linear system, this would be un-necessary work. Instead we wait to exit the `while` loop (i.e. we wait until the convergence criterion has been reached), and then solve the upper triangular linear system to compute the solution.
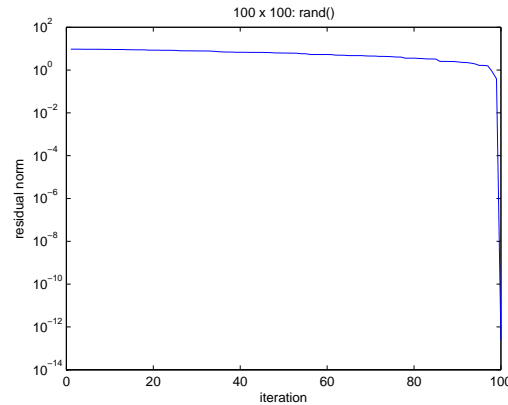
The last part of the code simply solves the linear system and returns the solution as the output of the function. Line 96 and 97 truncate the matrix, `HH` and the vector, `c`. These objects are then passed to the backslash operator (Gaussian Elimination) in line 100, to find the solution. Lastly, in line 105 we remove the newest column from the matrix of orthonormal basis vectors, then in line 107, we project the solution, `y` and add it to the initial guess to form our solution vector, which is the returned object.
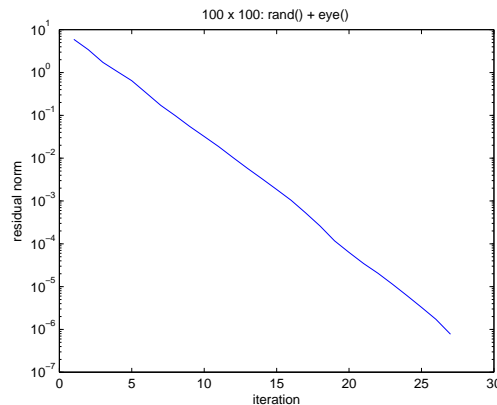
# 6 Testing the code

Whilst programming the matrix class and the linear solvers, I used test code to check that what was being produced was expected. The test code for the matrix and vector classes and for the Gaussian Elimination function is on the file named, 'use_matrix.cpp'. The tests I ran all produced the correct results and this file can be found in the appendix. The way I checked the Gaussian elimination solver was simply to multiply the result vector by the matrix and check that you ended up with the original RHS vector.

To test GMRES, I wrote three files called 'use_GMRES1.cpp', 'use_GMRES2.cpp' and 'use_GMRES3.cpp'. The first two are not really that interesting - they just compute the GMRES solution for small systems. 'use_GMRES3.cpp' builds a large matrix and fills it with uniformly distributed random numbers, using the `rand()` command. By modifying the GMRES routine ('GMRESout.cpp') to print the residual at each iteration to a data file, we can produce some convergence plots by loading the data into MATLAB.
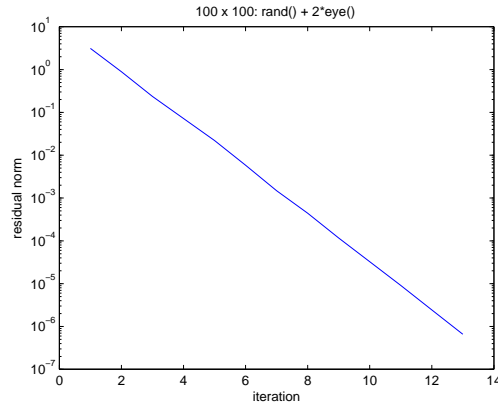
We will solve several $100 \times 100$ systems. Filling the entries with scaled random numbers makes for an ill-conditioned matrix and we thus expect convergence to take many iterations (100 or more). We can improve the conditioning of the matrix by adding multiples of the identity matrix, thus shifting the eigenvalues. This should have a dramatic effect upon convergence and if the GMRES method has been programmed correctly, we should be able to observe it on the convergence plots. First, the unshifted and badly conditioned matrix.
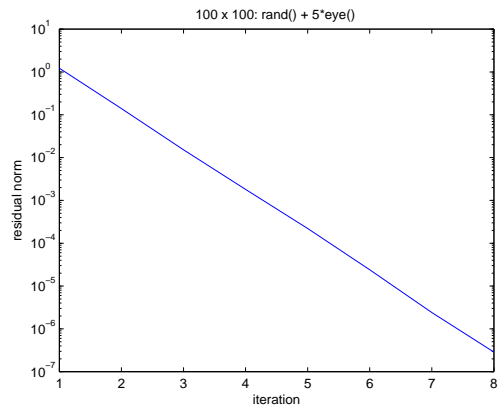


As expected, convergence (to our default tolerance of 1e-6) takes 100 iterations. Next, the same matrix with the identity matrix added on:

A dramatic improvement. Convergence now occurs in 27 iterations and is more uniform. Next, two times the identity added on:



100 x 100: rand() + 2*eye()

Convergence here occurs in 13 iterations. Lastly, we try adding 5 times the identity on:



100 x 100: rand() + 5*eye()

# 7    Extensions

There are many directions in which we could extend this work. Starting with GMRES, more features could be added, such as an optional maximum iteration parameter. It would also be a reasonable next step to code GMRES with restarts (we have programmed full GMRES here). Another obvious extension would be to program an effective error class to handle exceptions.

Aside form these natural extensions, there are a vast number of other linear solvers which could be programmed (CG, CGS, BICGStab etc). Each of these solvers often has a particular niche where they can perform better than the others, so combining all these methods into one bundle might be a reasonable challenge for a C++ programmer. Thinking even further ahead, it is certainly not impossible to conceive of designing a new GUI which would run these codes in a more user-friendly way.

# 8 Appendix - the codes in full

## 8.1 matrix.h

```
 1  #ifndef MATRIXDEF
 2  #define MATRIXDEF
 3
 4  ////////////////////////////////////////////////////////////////////////////
 5  /////////////     This header file defines the 'matrix' class    /////////////
 6  /////////////     All matrices created have 'double' entries     /////////////
 7  ////////////////////////////////////////////////////////////////////////////
 8
 9  #include <math.h>
10  #include <string>
11  #include <iostream>
12  #include <cassert>
13
14  using std::ostream;
15
16
17  class matrix
18  {
19
20  ////////////////////// Variables of the matrix class ////////////////////////
21  //private:
22  protected:
23     // Define dimensions of matrix
24        int rows, columns;
25        // Pointer to first entry of the vector of pointers
26        // (each of these pointers points to the first entry of each row)
27        double **xx;
28
29
30  ////////////////////// Constructors //////////////////////////////////////////
31
32  public:
33     // Overwritten default constructor
34        matrix();
35
36        // Creates matrix of given dimension
37        matrix(int no_of_rows,int no_of_columns);
38
39     // Overwritten copy constructor
40        matrix(const matrix& A);
41
42  ////////////////// Destructor /////////////////////////////////////////////
43
44        ¬matrix();
45
46  ////////////////// Binary Operators /////////////////////////////////////////
47
48        friend matrix operator +(const matrix& A, const matrix& B);
49        friend matrix operator -(const matrix& A, const matrix& B);
50
51        friend matrix operator *(const double& p, const matrix& A);
52        friend matrix operator *(const matrix& A, const double& p);
53        friend matrix operator *(const matrix& A, const matrix& B);
54
55        friend matrix operator /(const matrix& A, const double& p);
56        friend matrix operator /(const matrix& b, const matrix& A);
```

```cpp
57
58  ////////////////////// Unary operators ////////////////////////////////////
59      friend matrix operator +(const matrix& A);
60      friend matrix operator -(const matrix& A);
61
62      // Overload ¬ to mean transpose
63      friend matrix operator ¬(const matrix& A);
64
65  ////////////////////// Other operators /////////////////////////////////////
66
67    // Overloads the assignment operator, '='
68      matrix& operator =(const matrix &v);
69
70    // Overloads (), so A(i,j) returns the i,j entry a la MATLAB
71      double &operator() (int i, int j);
72
73      // Returns the row dimension of the matrix
74      friend int rows(matrix A);
75      // Returns the column dimension of the matrix
76      friend int columns(matrix A);
77
78  ///////////////////// Functions that are friends ////////////////////////
79
80      // Overloads the '<<' operator to allow easy printing of matrices
81      friend ostream& operator<<(ostream& output, const matrix& A);
82
83      // Create nxn Identity matrix
84      friend matrix eye(int size);
85
86    // computes an nxn permutation matrix which swaps rows i and j
87      friend matrix permute_r(int n, int i, int j);
88
89      // Locates largest number below the diagonal, for matrix, A
90      friend int find_pivot(matrix A, int column);
91
92      friend matrix resize(matrix A, int m, int n);
93  };
94  #endif
```

25

## 8.2 matrix.cpp

```cpp
1  #include "matrix.h"
2  #include "vector.h"
3  #include "error.h"
4
5  /////////////////// Constructors ///////////////////////////////////////////
6
7  // Constructor that overrides compiler generated default constructor
8  matrix::matrix()
9  {
10   // constructs an empty object with rows = columns = 0
11     rows = 0;
12     columns = 0;
13     xx = NULL;
14 }
15
16 // Constructor for basic matrix with specified dimensions
17 matrix::matrix(int no_of_rows, int no_of_columns)
18 {
19   // Matrix dimension fields
20     rows = no_of_rows;
21     columns = no_of_columns;
22
23   // A is an array of m pointers, each pointing to the
24   // first entry in the vector (of length, 'no_of_columns')
25     xx = new double *[no_of_rows];
26
27     // Allocate the memory for the entries of the matrix
28     for (int i=0; i<no_of_rows; i++)
29     {
30       // Creates 'no_of_rows' rows of length, 'no_of_columns'
31         xx[i] = new double[no_of_columns];
32     }
33
34     for (int i=0; i<no_of_rows; i++)
35     {
36         for (int j=0; j<no_of_columns; j++)
37         {
38           // Initialise the entries of the matrix to zero
39             xx[i][j]=0.0;
40         }
41     }
42 }
43
44 // Copy constructor - creates matrix with the same entries as input, A
45 matrix::matrix(const matrix& A)
46 {
47     rows = A.rows;                           // Matrix dimension fields
48     columns = A.columns;
49
50   // A is an array of m pointers, each pointing to the
51   // first entry in the vector (of length, 'A.columns')
52     xx = new double *[A.rows];
53
54   // Allocate the memory for the entries of the matrix
55     for (int i=0; i<A.rows; i++)
56     {
57       // Creates 'A.rows' rows of length, 'A.columns'
58         xx[i] = new double[A.columns];
```

```cpp
 59          }
 60
 61      for (int i=0; i < A.rows; i++)
 62      {
 63          for (int j=0; j < A.columns; j++)
 64          {
 65              // Copy across the entries from matrix A
 66              xx[i][j]=A.xx[i][j];
 67          }
 68      }
 69  }
 70
 71  // Destructor
 72  matrix::¬matrix()
 73  {
 74      if ( rows > 0 || columns > 0 )
 75      {
 76          for (int i=0; i < rows; i++)
 77          {
 78              delete[] xx[i];
 79          }
 80          delete[] xx;
 81      }
 82  }
 83
 84  /////////////////////// Binary Operators //////////////////////////////////////
 85
 86  // Overload the + operator to evaluate: A + B, where A and B are matrices
 87  matrix operator +(const matrix& A, const matrix& B)
 88  {
 89      int m = rows(A), n = columns(A), p = rows(B), q = columns(B);
 90
 91      if ( m != p || n !=q )
 92      {
 93          std::cout << "Error: Matrices of different dimensions";
 94          std::cout << "Returned first argument";
 95          return A;
 96      }
 97      else
 98      {
 99          matrix C(m,n);
100          for (int i=0; i<m; i++)
101          {
102              for (int j=0; j<n; j++)
103              {
104                  C.xx[i][j]=A.xx[i][j]+B.xx[i][j];
105              }
106          }
107          return C;
108      }
109  }
110
111
112
113  // Overload the − operator to evaluate: A − B, where A and B are matrices
114  matrix operator −(const matrix& A, const matrix& B)
115  {
116      int m,n,p,q;
117      m = rows(A);
118      n = columns(A);
119      p = rows(B);
120      q = columns(B);
```

```cpp
121
122      if ( m != p || n !=q )
123      {
124          std::cout << "Error: Matrices of different dimensions";
125          std::cout << "Returned first argument";
126          return A;
127      }
128      else
129      {
130          matrix C(m,n);
131          for (int i=0; i<m; i++)
132          {
133              for (int j=0; j<n; j++)
134              {
135                  C.xx[i][j]=A.xx[i][j]-B.xx[i][j];
136              }
137          }
138          return C;
139      }
140 }
141
142
143 // Definition of multiplication between a scalar, p and a matrix, A
144 matrix operator *(const double& p, const matrix& A)
145 {
146   // Create a matrix with the same dimensions as A
147     matrix B(A.rows,A.columns);
148
149     for (int i=0; i<A.rows; i++)
150     {
151         for (int j=0; j<A.columns; j++)
152         {
153             B.xx[i][j]= p * A.xx[i][j];     // Multiply each entry by p
154         }
155     }
156     return B;
157 }
158
159
160 // Definition of multiplication between a matrix, A and a scalar, p
161 matrix operator *(const matrix& A, const double& p)
162 {
163   // Create a matrix with the same dimensions as A
164     matrix B(A.rows,A.columns);
165
166     for (int i=0; i<A.rows; i++)
167     {
168         for (int j=0; j<A.columns; j++)
169         {
170             B.xx[i][j]= p * A.xx[i][j];     // Multiply each entry by p
171         }
172     }
173     return B;
174 }
175
176
177 // Definition of division of a matrix, A by a scalar, p i.e. A/p
178 matrix operator /(const matrix& A, const double& p)
179 {
180   // Create a matrix with the same dimensions as A
181     matrix B(A.rows,A.columns);
182
```

```
183     for (int i=0; i<A.rows; i++)
184     {
185         for (int j=0; j<A.columns; j++)
186         {
187             B.xx[i][j]= A.xx[i][j]/p;        // Divide each entry by p
188         }
189     }
190     return B;
191 }
192
193
194 // Define multiplication for matrices:
195
196 matrix operator *(const matrix& A, const matrix& B)
197 {
198   // Use assertion to check matrix dimensions are consistent
199     assert(A.columns==B.rows);
200
201   // Create a result matrix, C with the correct dimensions
202     matrix C(A.rows,B.columns);
203
204     double temp = 0;
205
206         // rows (m)
207         for (int i=0; i < A.rows ; i++)
208         {
209           // columns (q)
210             for (int j=0; j < B.columns ; j++)
211             {
212               // dot product step (n sums)
213                 for (int k=0; k < A.columns ; k++)
214                 {
215                     temp = temp + A.xx[i][k]*B.xx[k][j];
216                 }
217
218                 // Set the C matrix values
219                 C.xx[i][j]=temp;
220
221                 // reset temp
222                 temp = 0;
223             }
224         }
225
226     return C;
227 }
228
229 /////////////////////// Unary operators //////////////////////////////////////
230
231 matrix operator +(const matrix& A)      // Define the unary operator, '+'
232 {
233   // Create a temporary matrix with the same dimensions as A
234     matrix B(A.rows,A.columns);
235
236   // Set the entires of B to be the same as those in A
237     for (int i=0; i < A.rows; i++)
238     {
239         for (int j=0; j < A.columns; j++)
240         {
241             B.xx[i][j] = A.xx[i][j];
242         }
243     }
244     return B;
```

```cpp
245  }
246
247  matrix operator −(const matrix& A)        // Define the unary operator, '−'
248  {
249     // Create a temporary matrix with the same dimensions as A
250        matrix B(A.rows,A.columns);
251
252     // Set the entires of B to be the same as those in A
253        for (int i=0; i < A.rows; i++)
254        {
255            for (int j=0; j < A.columns; j++)
256            {
257                B.xx[i][j] = −A.xx[i][j];
258            }
259        }
260        return B;
261  }
262
263  // Overload the ¬ operator to mean transpose
264  matrix operator ¬(const matrix& A)
265  {
266     // Create a temporary matrix with reversed dimensions
267        matrix B(A.columns,A.rows);
268
269     // Set the entires of B to be the same as those in A
270        for (int i=0; i < A.columns; i++)
271        {
272            for (int j=0; j < A.rows; j++)
273            {
274                B.xx[i][j] = A.xx[j][i];
275            }
276        }
277        return B;
278  }
279
280  ///////////////////////// Other operators /////////////////////////
281
282  // Definition of matrix operator '='
283  // Operator returns a matrix equal to the RHS
284  matrix& matrix::operator =(const matrix &A)
285  {
286     // Destruct previous entries
287        for (int i=0; i < rows; i++)
288            {
289                delete[] xx[i];
290            }
291        delete[] xx;
292
293     // Assign new dimensions to be equal to that of the RHS
294        rows = A.rows;
295        columns = A.columns;
296
297     // Allocate the memory as in the constructor
298        xx = new double *[A.rows];
299
300        for (int i=0; i < A.rows ; i++)
301        {
302            xx[i] = new double[A.columns];
303        }
304
305     // Copy the values across from the RHS
306        for (int i=0; i < A.rows; i++)
```

```cpp
307          {
308              for (int j=0; j < A.columns ; j++)
309              {
310                // Set entries to be the same as the RHS matrix
311                  xx[i][j]=A.xx[i][j];
312              }
313          }
314          return *this;
315  }
316
317  // Allows reference to the entries of a matrix in the same way as MATLAB
318  // Can call or assign values.
319  double &matrix::operator() (int i, int j)
320  {
321      if (i < 1 || j < 1)
322      {
323      std::cout << "Error: One of your indices may have been too small \n\n";
324      }
325
326      else if (i > rows || j > columns)
327      {
328      std::cout << "Error: One of your indices may have been too large \n\n";
329      }
330      return xx[i-1][j-1];
331  }
332
333
334  //////////////////////// Function prototypes ////////////////////////////
335
336  int rows(matrix A);
337  int columns(matrix A);
338
339  //////////////////////// Function definitons ////////////////////////////
340
341  // Returns the private field, 'rows'
342  int rows(matrix A)
343  {
344      return A.rows;
345  }
346
347  // Returns the private field, 'columns'
348  int columns(matrix A)
349  {
350      return A.columns;
351  }
352
353  //////////////////////// Functions that are friends ///////////////////
354
355  // Overloads the '<<' operator to allow easy printing of matrices
356  ostream& operator << (ostream& output, const matrix &A)
357  {
358      for (int i=0; i<A.rows; i++)
359      {
360          for (int j=0; j < A.columns; j++)
361          {
362              output << " " << A.xx[i][j];
363          }
364          output << "\n";
365      }
366      output << "\n";
367      return output;
368  }
```

```
369
370  matrix eye(int size)
371  {
372    // Create a temporary matrix with the same dimensions as A
373      matrix temp_eye(size,size);
374
375    // Set the entries of B to be the same as those in A
376      for (int i=0; i < size; i++)
377      {
378              temp_eye.xx[i][i] = 1;
379      }
380      return temp_eye;
381  }
382
383  // Function that returns an nxn permutation matrix which swaps rows i and j
384  matrix permute_r(int n, int i, int j)
385  {
386    // Create nxn identity matrix
387      matrix I=eye(n);
388
389      // Zero the diagonal entries in the given rows
390      I(i,i)=0;
391      I(j,j)=0;
392
393    // Set the appropriate values to be 1
394      I(i,j)=1;
395      I(j,i)=1;
396
397      return I;
398  }
399
400  // Function that returns the row number of the largest
401  // sub-diagonal value of a given column
402
403  int find_pivot(matrix A, int column)
404  {
405    // Initialise maxval to be diagonal entry in column, 'column'
406      double maxval=fabs(A(column,column));
407
408      // Initialise rowval to be column
409      int rowval=column;
410
411          for (int i=column+1; i <= A.rows; i++)
412          {
413              if ( fabs(A(i,column)) > maxval)
414              {
415                  // Update maxval and rowval if bigger than previous maxval
416                  maxval = fabs(A(i,column));
417                  rowval = i;
418              }
419          }
420      return rowval;
421  }
422
423
424  // Function that returns an mxn matrix with entries that
425  // are the same as matrix A, where possible
426
427  matrix resize(matrix A, int m, int n)
428  {
429      int p,q;
430      matrix Mout(m,n);
```

```
431
432    // select lowest of each matrix dimension
433      if (m≤A.rows)
434      {
435          p=m;
436      }
437      else
438      {
439          p=A.rows;
440      }
441
442      if (n≤A.columns)
443      {
444          q=n;
445      }
446      else
447      {
448          q=A.columns;
449      }
450
451    // copy across relevant values
452      for (int i=1; i ≤ p; i++)
453      {
454          for (int j=1; j ≤ q; j++)
455          {
456              Mout(i,j) = A(i,j);
457          }
458      }
459
460      return Mout;
461  }
```

## 8.3 vector.h

```cpp
1  #ifndef VECTORDEF
2  #define VECTORDEF
3
4  #include "matrix.h"
5
6  // 'vector' inherits from 'matrix' in a public way
7  class vector: public matrix
8  {
9
10 public:
11
12 ///////////////////////// Constructors //////////////////////////////////////
13
14   // default constructor
15     vector();
16     // Constructor that takes 1 argument (size of the vector)
17     vector(int no_of_elements);
18
19 ////////////////////// Binary Operators //////////////////////////////////////
20
21     friend vector operator +(const vector& A, const vector& B);
22     friend vector operator -(const vector& A, const vector& B);
23
24     friend vector operator *(const double& p, const vector& A);
25     friend vector operator *(const vector& A, const double& p);
26
27     friend vector operator /(const vector& A, const double& p);
28
29
30 /////////////////////// Unary operators //////////////////////////////////////
31
32     friend vector operator +(const vector& A);
33     friend vector operator -(const vector& A);
34
35
36 /////////////////////// Other operators //////////////////////////////////////
37
38   // Overloads (), so x(i) returns the ith entry a la MATLAB
39     double &operator() (int i);
40
41     // Overloads the assignment operator, '=' for vector RHS
42     vector& operator =(const vector &v);
43
44 /////////////////////// Functions that are friends ///////////////////////
45
46     friend vector mat2vec(matrix A);
47
48   // Default call is norm(v) and returns 2-norm
49     friend double norm(vector v, int p=2);
50
51     friend vector GMRES(matrix A, matrix b, matrix x0, double tol=1e-6);
52
53     friend vector GMRESout(matrix A, matrix b, matrix x0, double tol=1e-6);
54
55     //friend vector resize(vector v, int m);
56
57 };
58 #endif
```

## 8.4   vector.cpp

```cpp
1   #include "vector.h"                    // include the header file
2
3
4   //////////////////////// Constructors //////////////////////////////////////
5
6   // in the class, 'vector' there is a constructor of the same name
7   vector::vector()
8
9   // runs the default matrix constructor
10  : matrix()
11  {
12
13  }
14
15  vector::vector(int no_of_elements)
16
17  : matrix(no_of_elements,1)
18  {
19
20  }
21
22  ///////////////////// Other operators ///////////////////////////////////
23
24  // Overloads (), so x(i) returns the ith entry a la MATLAB
25  double &vector::operator() (int i)
26  {                                        // Can call or assign values.
27      if (i < 1)
28      {
29          std::cout << "Error: Your index may be too small \n\n";
30      }
31
32      else if (i > rows)
33      {
34          std::cout << "Error: Your index may be too large \n\n";
35      }
36      return xx[i−1][0];
37  }
38
39  // Operator returns a matrix equal to the RHS
40  vector& vector::operator =(const vector &v)
41  {
42    // Destruct previous entries
43      for (int i=0; i < rows; i++)
44          {
45              delete[] xx[i];
46          }
47      delete[] xx;
48
49    // Assign new dimensions to be equal to that of the RHS
50      rows = v.rows;
51      columns = v.columns;
52
53    // Allocate the memory as in the constructor
54      xx = new double *[v.rows];
55
56      for (int i=0; i < v.rows ; i++)
57      {
58          xx[i] = new double[v.columns];
```

```
59        }
60
61    // Copy the values across from the RHS
62        for (int i=0; i < v.rows; i++)
63        {
64            for (int j=0; j < v.columns ; j++)
65            {
66              // Set entries to be the same as the RHS matrix
67                xx[i][j]=v.xx[i][j];
68            }
69        }
70        return *this;
71
72 }
73 ////////////////// Binary Operators /////////////////////////////////////////
74
75 vector operator +(const vector& A, const vector& B)
76 {
77        int m,n;
78        m = rows(A);
79        n = rows(B);
80
81        if ( m != n )
82        {
83            std::cout << "Error: Matrices of different dimensions.";
84            std::cout << " Returned first argument";
85            return A;
86        }
87        else
88        {
89            vector v(m);
90            for (int i=0; i<m; i++)
91            {
92                v(i+1) = A.xx[i][0]+B.xx[i][0];
93            }
94            return v;
95        }
96 }
97
98 vector operator −(const vector& A, const vector& B)
99 {
100       int m,n;
101       m = rows(A);
102       n = rows(B);
103
104       if ( m != n )
105       {
106           std::cout << "Error: Matrices of different dimensions.";
107           std::cout << " Returned first argument";
108           return A;
109       }
110       else
111       {
112           vector v(m);
113           for (int i=0; i<m; i++)
114           {
115               v(i+1) = A.xx[i][0]−B.xx[i][0];
116           }
117           return v;
118       }
119 }
120
```

```
121  //
122  vector operator *(const double& p, const vector& A)
123  {
124          int m = rows(A);
125          vector v(m);
126          for (int i=0; i<m; i++)
127          {
128              v(i+1) = p*A.xx[i][0];
129          }
130          return v;
131  }
132  //
133  vector operator *(const vector& A, const double& p)
134  {
135          int m = rows(A);
136          vector v(m);
137          for (int i=0; i<m; i++)
138          {
139              v(i+1) = p*A.xx[i][0];
140          }
141          return v;
142  }
143  //
144  vector operator /(const vector& A, const double& p)
145  {
146          int m = rows(A);
147          vector v(m);
148          for (int i=0; i<m; i++)
149          {
150              v(i+1) = A.xx[i][0]/p;
151          }
152          return v;
153  }
154  //
155
156  ///////////////////////// Unary operators ////////////////////////////////
157
158  vector operator +(const vector& A)
159  {
160          int m = rows(A);
161          vector v(m);
162          for (int i=0; i<m; i++)
163          {
164              v(i+1) = A.xx[i][0];
165          }
166          return v;
167  }
168  //
169  vector operator −(const vector& A)
170  {
171          int m = rows(A);
172          vector v(m);
173          for (int i=0; i<m; i++)
174          {
175              v(i+1) = −A.xx[i][0];
176          }
177          return v;
178  }
179
180  ///////////////////////// Functions that are friends ////////////////////////
181
182  // Function that returns the first column of a matrix, A as a vector
```

```
183
184  vector mat2vec(matrix A)
185  {
186     // create vector with same no of rows as A and 1 column
187        vector v(rows(A));
188
189        for (int i=1; i <= rows(A); i++)
190        {
191            v(i)=A(i,1);                             // copy only first column
192        }
193
194        return v;
195  }
196
197  double norm(vector v, int p)
198  {
199     // define variables and initialise sum
200        double temp, value, sum = 0.0;
201
202        for (int i=1; i <= rows(v); i++)
203        {
204            // floating point absolute value
205            temp = fabs(v(i));
206            sum += pow(temp, p);
207        }
208
209        value = pow(sum, 1.0/((double)(p)));
210
211        return value;
212  }
```

## 8.5   backslash.cpp

```cpp
#include "matrix.h"
#include "vector.h"

// Definition of division of a vector, b by a matrix, A i.e. y=b/A
matrix operator /(const matrix& b, const matrix& A)
{
    int n = A.rows;

  // Create empty matrices, P & L
    matrix P, L;

    // Create and intialise U & Atemp
    matrix Utemp = eye(n);
    matrix Atemp=A;

    //std::cout << U << "\n\n";

    for (int j=1; j < n ; j++)
    {

        //std::cout << "Need to permute row " << j << " with row ";
    //std::cout << find_pivot(Atemp,j) << "\n\n";

    // Create appropriate permutation matrix, P
        P = permute_r(n,find_pivot(Atemp,j),j);

        Utemp = P*Utemp;                                 // Update U & Atemp

        Atemp = Utemp*A;

        //std::cout << "Permute rows \n\n" << Atemp;

        L = eye(n);

        for (int i=j+1; i <= n ; i++)
        {
          // Check for division by zero
            assert(fabs(Atemp(j,j))>1.0e-015);

            // Compute multiplier and store in sub-diagonal entry of L
            L(i,j)= -Atemp(i,j)/Atemp(j,j);
        }

        Utemp = L*Utemp;

        Atemp = Utemp*A;

        //std::cout << "Eliminate sub-diagonal entries \n\n" << Atemp;

    }

// Now loop through and set to zero any values which are almost zero

    for (int j=1; j < n ; j++)
    {
        for (int i=j+1; i <= n ; i++)
        {
            if (fabs(Atemp(i,j)) < 5.0e-016)
```

```
59                {
60                    Atemp(i,j)=0;
61                }
62            }
63        }
64
65        //std::cout << "The matrix U = Utemp*A is then: \n\n" << Atemp;
66
67        // So, to solve Ax=b, we do: (Utemp*A)x=Utemp*b i.e.
68        // Set U=Utemp*A=Atemp, compute y=Utemp*b and
69        // solve Ux=y (upper triangular system -> back subs)
70
71        matrix U = Utemp*A;             //Atemp; gives the same result
72        matrix y = Utemp*b;
73
74        //std::cout << "The RHS is then Utemp*b: \n\n" << y ;
75
76        matrix x(n,1);                  // Create result vector
77
78        // Solve Ux=y by back substitution:
79
80     // Compute last entry of vector x (first step in back subs)
81        x(n,1)=y(n,1)/U(n,n);
82
83        double temp = 0;                // Initialise temp
84
85        for (int i=n-1; i>=1; i--)
86        {
87            temp = y(i,1);
88            for (int j=n; j>i; j--)
89            {
90                temp = temp - U(i,j)*x(j,1);
91            }
92            x(i,1)=temp/U(i,i);
93        }
94
95        return x;
96 }
```

## 8.6  GMRES.cpp

```cpp
1   #include "matrix.h"
2   #include "vector.h"
3
4   // GMRES
5   vector GMRES(matrix A, matrix b, matrix x0, double tol)
6   {
7       //double tol=1e-6;
8
9     // determine initial residual, r0 in vector form
10      vector r0 = mat2vec(b - A*x0);
11
12      //std::cout << "initial residual vector, r0 = b-A*x0 : \n\n" << r0;
13
14    // need this in least square part later
15      double normr0 = norm(r0);
16
17      // initialise to enter while loop
18      double residual=1.0;
19
20      // intialise vector v
21      vector v= r0/normr0;
22
23      //std::cout << "initial vector v = r0/||ro|| : \n\n" << v;
24
25    // Arnoldi/GMRES step index
26      int k=1;
27
28      // Declare Givens rotation matrix, initialise Jtotal;
29      matrix J, Jtotal;
30      Jtotal=eye(2);
31
32    // intialise H, declare tempMat, V, w
33      matrix H(1,1), Htemp, HH, bb(1,1), c, cc;
34      matrix tempMat, V, Vold, hNewCol;
35      vector w, vj(rows(v));
36
37      bb(1,1)=normr0;
38
39    // initialise matrix V (matrix of orthogonal basis vectors)
40      V=v;
41
42      while (residual>tol)
43      {
44          //std::cout<< " ----------------------------------------------\n\n";
45
46      // update Vold (used for checking Arnoldi later)
47          Vold=V;
48
49          H=resize(H,k+1,k);
50
51      // Arnoldi steps (using Gram-Schmidt process)
52          w = mat2vec(A*v);
53          //std::cout<< "(k = " << k <<") : vector w=Av : \n\n" << w;
54
55              for (int j=1; j≤k; j++)
56              {
57                  for (int i=1; i≤rows(V); i++)
58                  {
```

41

```cpp
59                    // set the vector vj to be jth column of V
60                        vj(i)=V(i,j);
61                    }
62
63                    tempMat = ¬vj*w;
64
65                    // these two lines calculate the inner product
66                    H(j,k)= tempMat(1,1);
67                    //std::cout<< "H("<<j<<","<<k<<")= "<<H(j,k)<<"\n\n";
68
69                    w = w − H(j,k)*vj;
70                    //std::cout<< "Gramm−Schmidt update of vector w: \n\n" << w;
71                }
72
73          H(k+1,k)=norm(w);
74          //std::cout<< "H(" << k+1 << "," << k << ")= " << H(k+1,k) << "\n\n";
75
76          v=w/H(k+1,k);
77          //std::cout<< "(k = " << k <<") :new vector v: \n\n" << v;
78
79      // add one more column to matrix V
80          V=resize(V,rows(V),k+1);
81
82          for (int i=1; i≤rows(V); i++)
83          {
84            // copy entries of v to new column of V
85              V(i,k+1)=v(i);
86          }
87
88          //std::cout<< "(k = " << k << ") :latest matrix, V: \n\n" << V;
89
90          //std::cout << "(k = " << k <<") :latest matrix, H: \n\n" << H;
91
92          //std::cout << "check:    AV[k] = V[k+1]H: \n\n" << A*Vold << V*H;
93
94      ///////////////////////////// Least squares step /////////////////////////
95
96          if (k==1)
97          {
98            // First pass through, Htemp=H
99              Htemp=H;
100         }
101         else
102         {
103           // for subsequent passes, Htemp=Jtotal*H
104             Jtotal=resize(Jtotal,k+1,k+1);
105             Jtotal(k+1,k+1)=1;
106             Htemp=Jtotal*H;
107         }
108
109     // Form next Givens rotation matrix
110         J = eye(k−1);
111         J = resize(J,k+1,k+1);
112
113     J(k,k)=Htemp(k,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
114     J(k,k+1)=Htemp(k+1,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
115     J(k+1,k)=−Htemp(k+1,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
116     J(k+1,k+1)=Htemp(k,k)/pow(pow(Htemp(k,k),2)+pow(Htemp(k+1,k),2),0.5);
117
118         //std::cout<< "J: \n\n" << J;
119
120     // combine together with previous Givens rotations
```

```
121        Jtotal=J*Jtotal;
122
123        //std::cout<< "Check orthogonality of Jtotal \n\n" << ¬Jtotal*Jtotal;
124
125        HH=Jtotal*H;
126
127        for (int i=1; i≤k+1; i++)
128        {
129            for (int j=1; j≤k; j++)
130            {
131              // set all 'small' values to zero
132                if (fabs(HH(i,j))<1e−15)
133                {
134                    HH(i,j)=0;
135                }
136            }
137        }
138
139        //std::cout<< "Check Jtotal*H is upper triangular: \n\n" << HH;
140
141        bb=resize(bb,k+1,1);
142
143        //std::cout<< "bb: \n\n" << bb;
144
145        c=Jtotal*bb;
146
147        //std::cout<< "c=J*bb: \n\n" << c;
148
149        residual=fabs(c(k+1,1));
150
151        //std::cout<< k << "th residual: \n\n" << residual << "\n\n";
152
153        k++;
154    }
155
156    std::cout<< "GMRES iteration converged in " << k−1 << " steps\n\n";
157
158  // Extract upper triangular square matrix
159    HH=resize(HH,rows(HH)−1,columns(HH));
160
161    //std::cout<< "HH: \n\n" << HH;
162
163    cc=resize(c,rows(HH),1);
164
165    //std::cout<< "cc: \n\n" << cc;
166
167    matrix yy = cc/HH;                                // solve linear system
168
169    vector y = mat2vec(yy);
170
171    //std::cout<< "y: \n\n" << y;
172
173  // chop the newest column off of matrix V
174    V=resize(V,rows(V),columns(V)−1);
175
176    vector x= mat2vec(x0+V*y);
177
178    return x;
179 }
```

## 8.7 use_matrix.cpp

```cpp
1   #include "matrix.h"
2
3   #include "vector.h"
4   #include <stdlib.h>
5
6   int use_matrix()
7   {
8       int m = 4, n = 4;
9
10      matrix A(m,n);
11
12      std::cout << "Matrix A has " << rows(A) << " rows" <<"\n";
13      std::cout << "Matrix A has " << columns(A) << " columns" <<"\n\n";
14
15      std::cout << "The newly created matrix, A looks like: \n\n";
16
17      std::cout << A ;
18
19      std::cout << "Now, we will assign some values to the entries: \n\n";
20
21      for (int i=1; i≤m; i++)
22      {
23          for (int j=1; j≤n; j++)
24          {
25              A(i,j)=i+j;
26          }
27      }
28
29      std::cout << A ;
30
31      std::cout << "The (1,1) entry of matrix A is: " << A(1,1) << "\n\n";
32
33      std::cout << "matrix, B has 5 times the entries of A: \n\n";
34
35      matrix B(m,n);
36
37      for (int i=1; i≤m; i++)
38      {
39          for (int j=1; j≤n; j++)
40          {
41              B(i,j)=5*(i+j);
42          }
43      }
44
45      std::cout << B ;
46
47      std::cout << "Add together the two matrices (A+B): \n\n";
48
49      std::cout << A+B;
50
51      std::cout << "Subtract the matrices (A—B): \n\n";
52
53      std::cout << A—B;
54
55      std::cout << "Create matrix C = A + A + A  using the '=', \n\n";
56
57      matrix C(m,n);
58      C = A + A + A;
```

```
59
60      std::cout << C;

62      std::cout << "Create matrix D such that D = +C \n\n";

64      matrix D(m,n);
65      D = +C;
66      std::cout << D;

68      std::cout << "Create matrix E such that E = -C \n\n";

70      matrix E(m,n);
71      E = -C;
72      std::cout << E;

74      std::cout << "Try copying the above matrix: \n\n";

76      matrix F(m,n);
77      F = matrix(E);
78      std::cout << F;

80      std::cout << "G = A is automatically sized: \n\n";

82      matrix G = A;

84      std::cout << G;

86      std::cout << "Create matrix H = 6*G: \n\n";

88      matrix H = 6 * G;

90      std::cout << H;

92      std::cout << "Create matrix J = H*0.5: \n\n";

94      matrix J = H * 0.5;

96      std::cout << J;

98      std::cout << "Overwrite matrix B = J/10: \n\n";

100     B = J/10;

102     std::cout << B;

104     //std::cout << "Now consider multiplying two Matrices together \n";
105     //std::cout << "Define L=A and M= \n\n";

107     //std::cout << A*B;

109     std::cout << "Create a 'vector', x: \n\n";

111     matrix x(n,1);

113     for (int i=1; i≤n; i++)
114     {
115         x(i,1)=i;
116     }

118     std::cout << x;

120     std::cout << "Multiply the Matrix A by x (Ax)  \n\n";
```

```
121
122      std::cout << A*x ;
123
124      ////////////////////////// Vector Stuff //////////////////////////////
125
126      std::cout << "Create an empty vector using default constructor \n\n";
127
128      vector a;
129
130      std::cout << "Vector a has " << rows(a) << " rows" <<"\n";
131      std::cout << "Vector a has " << columns(a) << " columns" <<"\n\n";
132
133      std::cout << "Create a vector of size " << n << "\n\n";
134
135      vector b(n);
136
137      std::cout << "Vector b has " << rows(b) << " rows" <<"\n";
138      std::cout << "Vector b has " << columns(b) << " columns" <<"\n\n";
139
140      std::cout << b;
141
142      std::cout << "Put some values into the entries: \n\n";
143
144      for (int i=1; i≤n; i++)
145      {
146          //b(i,1)=i;
147          b(i)=i;
148      }
149
150      std::cout << b;
151
152      std::cout << "Multiply matrix A by vector b \n\n";
153
154      std::cout << A*b;
155
156      std::cout << "Create a 7 x 7 identity matrix \n\n";
157
158      B = eye(7);
159
160      std::cout << B;
161
162      matrix T(4,4); A=T;
163
164      std::cout << "Perform GE w. PP on the following 4x4 matrix: \n\n";
165
166      A(1,1)=2; A(1,2) = 1; A(1,3) = 1; A(1,4) = 0;
167      A(2,1)=4; A(2,2) = 3; A(2,3) = 3; A(2,4) = 1;
168      A(3,1)=8; A(3,2) = 7; A(3,3) = 9; A(3,4) = 5;
169      A(4,1)=6; A(4,2) = 7; A(4,3) = 9; A(4,4) = 8;
170
171      std::cout << "Matrix A is: \n\n" << A;
172
173      std::cout << "Vector b is: \n\n" << b << "Solve x=b/A: \n\n";
174
175      x=b/A;
176
177      std::cout << "The solution to the problem, x is: \n\n" << x;
178
179      std::cout << "And as a check, multiply A*x: \n\n";
180
181      std::cout << A*x;
182
```

```cpp
183        std::cout << "Try another example: \n\n";
184
185        matrix AA(3,3);
186
187        AA(1,1)=3; AA(1,2) = 17; AA(1,3) = 10;
188        AA(2,1)=2; AA(2,2) = 4; AA(2,3) = −2;
189        AA(3,1)=6; AA(3,2) = 18; AA(3,3) = −12;
190
191        vector c(3);
192        c(1)=1; c(2)=2; c(3) =3;
193
194        std::cout << "Matrix AA is: \n\n" << AA << "Vector c is \n\n" << c;
195
196        matrix y=c/AA;
197
198        std::cout << "The solution to the problem, y is: \n\n" << y;
199
200        std::cout << "And as a check, multiply A*y: \n\n";
201
202        std::cout << AA*y;
203
204        /*x=matrix(2,2);
205        x(1,1)=1;x(1,2)=2;x(2,1)=3;x(2,2)=4;
206        x=2*x;*/
207
208        std::cout<<A;
209
210        A=resize(A,10,3);
211
212        std::cout<<AA;
213
214        std::cout<<¬AA;
215
216        vector d(10);
217        for (int i=1; i≤10; i++)
218        {
219            d(i)=i;
220        }
221
222        std::cout<<d;
223
224        std::cout<<¬d*d;
225
226    exit(0);
227 }
```

## 8.8 useGMRES.cpp

```cpp
1  #include "matrix.h"
2  #include "vector.h"
3  //#include "error.h"
4  #include <stdlib.h>
5
6  int useGMRES()
7  {
8
9    matrix A(5,5);
10   //A(1,5)=1;
11   //A(2,1)=1;
12   //A(3,2)=1;
13   //A(4,3)=1;
14   //A(5,4)=1;
15
16 A(1,1)=0.8780;A(1,2)=0.8316;A(1,3)=0.2663;A(1,4)=0.9787;A(1,5)=0.0239;
17 A(2,1)=0.1159;A(2,2)=0.2926;A(2,3)=0.2626;A(2,4)=0.7914;A(2,5)=0.2085;
18 A(3,1)=0.9857;A(3,2)=0.5109;A(3,3)=0.5826;A(3,4)=0.2115;A(3,5)=0.2943;
19 A(4,1)=0.8573;A(4,2)=0.7512;A(4,3)=0.4431;A(4,4)=0.9486;A(4,5)=0.3660;
20 A(5,1)=0.4416;A(5,2)=0.3803;A(5,3)=0.4465;A(5,4)=0.0586;A(5,5)=0.8501;
21
22   vector b(5);
23   //b(1)=1;
24
25   b(1)=1;b(2)=1;b(3)=1;b(4)=1;b(5)=1;
26
27   vector x0(5);
28
29   std::cout << "matrix A: \n\n" << A;
30   std::cout << "vector b: \n\n" << b;
31   std::cout << "initial guess vector, x0: \n\n" << x0;
32
33   vector x = GMRES(A,b,x0);
34
35   std::cout << "GMRES solution, x is: \n\n" << x;
36
37   std::cout << "Check: Ax = \n\n" << A*x;
38
39   std::cout << "Backslash solution b/A \n\n" << b/A;
40
41   srand(7);
42
43   std::cout << "Random number 1:   " << - 1000 + rand() 2000  << " \n\n";
44   std::cout << "Random number 2:   " << - 1000 + rand() 2000  << " \n\n";
45   std::cout << "Random number 3:   " << - 1000 + rand() 2000  << " \n\n";
46
47   exit(0);
48 }
```

## 8.9 useGMRES2.cpp

```cpp
#include "matrix.h"
#include "vector.h"
//#include "error.h"
#include <stdlib.h>

int useGMRES2()
{
  matrix A(4,4);
  //A(1,5)=1;
  //A(2,1)=1;
  //A(3,2)=1;
  //A(4,3)=1;
  //A(5,4)=1;

A(1,1)=0.8780;A(1,2)=0.8316;A(1,3)=0.2663;A(1,4)=0.9787;//A(1,5)=0.0239;
A(2,1)=0.1159;A(2,2)=0.2926;A(2,3)=0.2626;A(2,4)=0.7914;//A(2,5)=0.2085;
A(3,1)=0.9857;A(3,2)=0.5109;A(3,3)=0.5826;A(3,4)=0.2115;//A(3,5)=0.2943;
A(4,1)=0.8573;A(4,2)=0.7512;A(4,3)=0.4431;A(4,4)=0.9486;//A(4,5)=0.3660;
//A(5,1)=0.4416;A(5,2)=0.3803;A(5,3)=0.4465;A(5,4)=0.0586;A(5,5)=0.8501;

  vector b(4);
  //b(1)=1;

  b(1)=1;b(2)=1;b(3)=1;b(4)=1;//b(5)=1;

  vector x0(4);

  std::cout << "matrix A: \n\n" << A;
  std::cout << "vector b: \n\n" << b;
  std::cout << "initial guess vector, x0: \n\n" << x0;

  vector x = GMRES(A,b,x0);

  std::cout << "GMRES solution, x is: \n\n" << x;

  std::cout << "Check: Ax = \n\n" << A*x;

  std::cout << "Backslash solution b/A \n\n" << b/A;

  exit(0);
}
```

## 8.10 useGMRES3.cpp

```cpp
1  #include "matrix.h"
2  #include "vector.h"
3  //#include "error.h"
4  #include <stdlib.h>
5
6  int useGMRES3()
7  {
8    int n=100;
9
10   matrix A(n,n);
11
12   srand(1);
13
14   for (int i=1; i≤n; i++)
15   {
16     for (int j=1; j≤n; j++)
17     {
18         A(i,j) = − 1000 + rand() % 2000;
19     }
20   }
21
22   A=A/1000;
23
24   A=A/pow(n,0.5);
25
26   A=A+5*eye(n);
27
28   vector b(n);
29
30   for (int i=1; i≤n; i++)
31   {
32         b(i) = 1.0;
33   }
34
35   vector x0(n);
36
37   //std::cout << "matrix A: \n\n" << A;
38   //std::cout << "vector b: \n\n" << b;
39   //std::cout << "initial guess vector, x0: \n\n" << x0;
40
41   vector x = GMRESout(A,b,x0);
42
43   std::cout << "GMRES solution, x is: \n\n" << x;
44
45   std::cout << "Check: Ax = \n\n" << A*x;
46
47   //std::cout << "Backslash solution b/A \n\n" << b/A;
48
49   exit(0);
50 }
```

# References

[1]  Lloyd N. Trefethen & David Bau III
     Numerical Linear Algebra
     SIAM 1997

[2]  Youcef Saad & Martin H. Schultz GMRES: A generalised minimal residual algorithm
     for solving nonsymmetric linear systems
     SIAM Journal of Scientific Computing., Vol. 7, No. 3, July 1986

[3]  Joe Pitt-Francis
     Lecture Notes
     C++ for Scientifc Computing Course 2009
     Computing Laboratory, Oxford University