

# The Conjugate Gradient Method

Will Woolfenden

November 12, 2023

## 1 Introduction

The Conjugate Gradient (CG) method is a method for efficiently solving a system of linear equations. It is efficient because time complexity is reduced from other methods, such as gradient descent or Newton's method. Linear systems with sparse matrices are common when we apply finite difference methods to solve boundary value problems. In this report we look at solving these linear systems, providing finite difference approximations to Poisson's equation in 1D  $u''(s) = g(s)$  and in 2D  $\nabla^2 u(s, t) = 1$ .

## 2 Methods, Results and Verification

### 2.1 Conjugate Gradient

In CG, we start with a guess  $\mathbf{x}_0$  to solve  $\mathbf{Ax} = \mathbf{b}$ . At every step, we compute  $\mathbf{r}_i = \mathbf{Ax}_i - \mathbf{b}$  with stopping criterion  $\|\mathbf{r}\| < \varepsilon$ . Our search vectors  $\mathbf{p}_i$  are chosen such that they obey  $\mathbf{p}_j^T \mathbf{Ap}_i$  for  $i \neq j$ , with  $\mathbf{p}_0 = \mathbf{r}_0$ . We perform this with a Gram-Schmidt orthogonalisation<sup>1</sup> with the inner product  $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{A}} = \mathbf{u}^T \mathbf{Av}$ . Each iteration performs

$$\begin{aligned}\mathbf{x}_{i+1} &= \mathbf{x}_i + \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{Ap}_i} \mathbf{p}_i \\ \mathbf{r}_{i+1} &= \mathbf{r}_i - \frac{\mathbf{r}_i^T \mathbf{r}_i}{\mathbf{p}_i^T \mathbf{Ap}_i} \mathbf{Ap}_i \\ \mathbf{p}_{i+1} &= \mathbf{r}_{i+1} + \frac{\mathbf{r}_{i+1}^T \mathbf{r}_{i+1}}{\mathbf{r}_i^T \mathbf{r}_i} \mathbf{p}_i.\end{aligned}$$

The residual vector  $\mathbf{r}_i$  describes the distance of  $\mathbf{x}_i$  to the solution  $\mathbf{x}^*$  which satisfies  $\mathbf{Ax}^* = \mathbf{b}$ . The vectors  $\mathbf{p}_i$  are the search directions of the iteration. For every iteration, we compute the distance we need to travel along each search direction, by intuitively realising that at a point  $\mathbf{x}_i$ , the vector from  $\mathbf{x}_i$  to the optimiser  $\mathbf{x}^*$  is  $A$ -orthogonal to all the search directions we have already traversed.

### 2.2 Finite Difference Methods: Poisson's Equation

We want to form a finite difference approximation to Poisson's equation in 1D  $u''(s) = g(s)$ , subject to boundary conditions  $u(0) = u(1) = 0$ . We discretise  $s$  into  $n + 2$  points<sup>2</sup> in the interval  $[0, 1]$

---

<sup>1</sup>this requires  $\mathbf{A}$  to be  $n \times n$  with all vectors being  $n$ -vectors.

<sup>2</sup> $n$  points excluding 0 and 1.

where  $s_k = kh$ ,  $k = 0, 1, 2, \dots, n + 1$  and  $h = 1/(n + 1)$ . The finite difference approximation to the second derivative is

$$u''(s) \approx \frac{u(s - h) - 2u(s) + u(s + h)}{h^2} \quad (1)$$

which can be seen by the Taylor series expansions of the  $u(s)$  terms. Write  $u_k = u(s_k)$ ,  $g_k = g(s_k)$ . Boundary conditions are  $u_0 = u_n = 0$ . We obtain the linear system

$$\begin{bmatrix} -2 & 1 & & & & & \\ 1 & -2 & 1 & & & & \\ & & 1 & \ddots & \ddots & & \\ & & & \ddots & \ddots & 1 & \\ & & & & 1 & -2 & 1 \\ & & & & & 1 & -2 \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ \vdots \\ u_{n-1} \\ u_n \end{pmatrix} = -h^2 \begin{pmatrix} g_1 \\ g_2 \\ \vdots \\ \vdots \\ g_{n-1} \\ g_n \end{pmatrix}.$$

This is a linear system of the form  $\mathbf{Ax} = \mathbf{b}$  which we can solve using the CG method. Consider the case where  $g$  is the constant function  $g(s) = -1$ , and where  $\mathbf{A}$  is  $5 \times 5$ . The system in full, changing sign, is

$$\begin{bmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix} \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \\ u_5 \end{pmatrix} = \frac{1}{(n + 1)^2} \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{pmatrix}.$$

Our code produces `tridiagonal_convergence_0611.txt`:

Time t	L2 norm	LInf norm
0	2.22448	1.97222
5.4392e-05	0.588795	0.397175
6.394e-05	0.233913	0.188237
7.0762e-05	0.10518	0.0963511
7.7545e-05	0.00941914	0.00515077
8.3386e-05	2.94053e-16	1.73472e-16

which claims to solve the system in exactly 5 iterations. This is appropriate for CG, which is designed to converge in no more iterations than the size of the system [3]. Our program prints the solution  $x = (0.0694444, 0.111111, 0.125, 0.111111, 0.0694444)$ . Using the solver in `Octave` verifies this computation, shown in B.1, hence it is the true solution of the linear system.

The ODE we are solving is  $u''(s) = -1$ . By direct integration, the general solution is  $u(s) = Cs^2 + Ds$  for constants  $C, D$ . Plugging in boundary conditions, we obtain the particular solution  $u(s) = -s(s - 1)/2$  in closed form. See Figure 1 for a comparison.

### 2.3 Perturbed Tridiagonal

We perturb the system such that the main diagonal contains  $2 + \alpha$  instead of 2, for some positive real constant  $\alpha$ . In this case, we impose  $\mathbf{b}$  to be the vector with all entries  $b_i = 2.5$ , equivalently all  $g_i = -5(n + 1)^2/2$ . Figure 2 illustrates the convergence of these tridiagonal systems. We notice that, as the entries on the main diagonal grow, the effectiveness of the iterations increases. The convergence of the CG method is  $\mathcal{O}(m\sqrt{\kappa})$ , where  $m$  is the number of non-zero entries in the matrix and  $\kappa$  is the condition number for any choice of norm [4]. As the entries on the diagonal grow,  $\mathbf{A}$

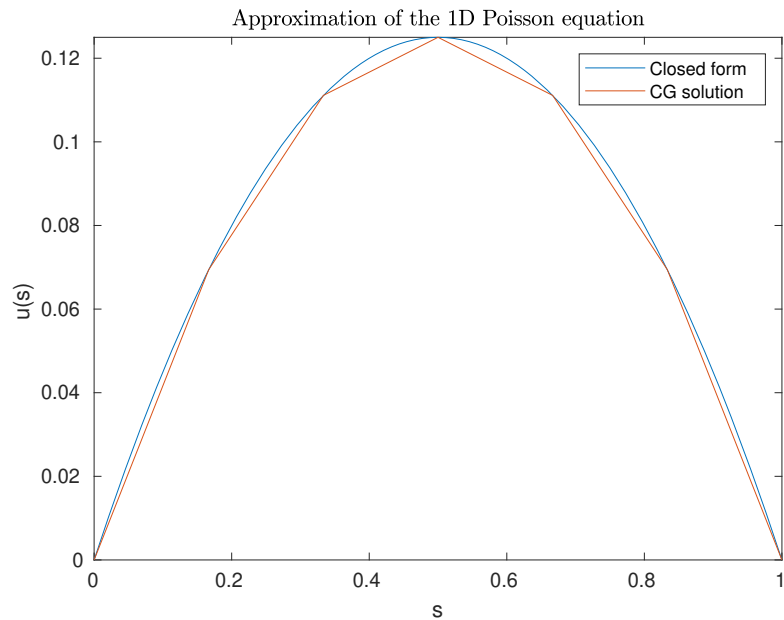


Figure 1: Comparison of the closed form solution to the ODE, using `fplot()`, and the  $n = 5$  finite difference approximation from solving the system with CG.

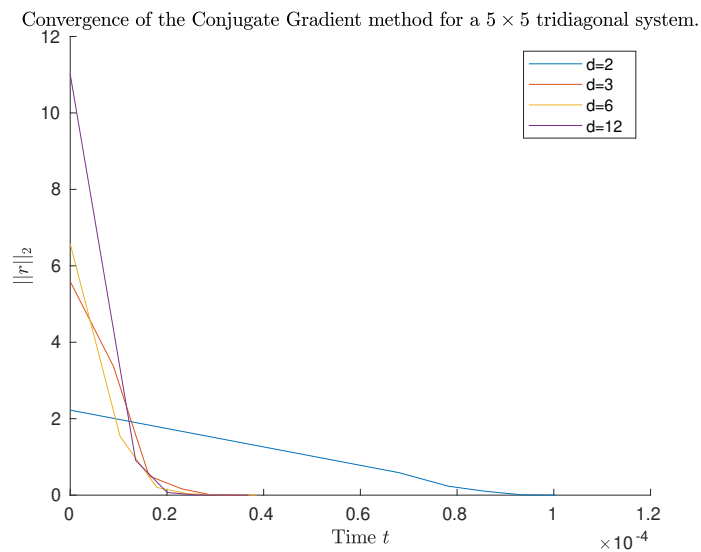


Figure 2: Convergence graphs for different size tridiagonal systems. The solutions converge more rapidly for larger diagonal entries. The matrices  $\mathbf{A}$  are the same as earlier except for  $2 + \alpha$  on the diagonal, and the vectors  $\mathbf{b}$  are  $b_k = 2.5$  for all  $k$ .

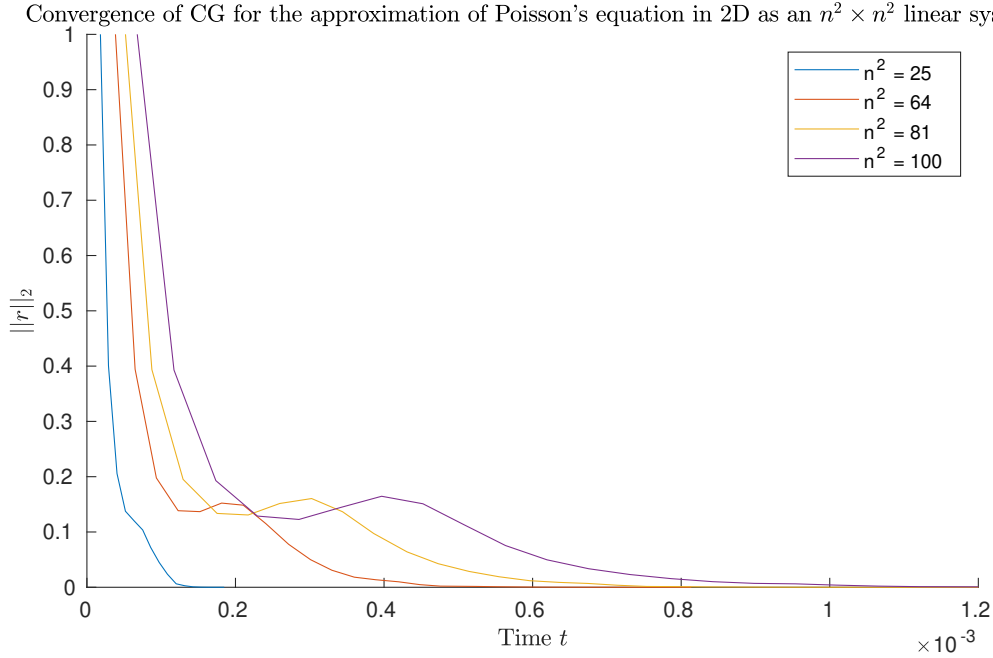


Figure 3: Convergence graphs for the CG method applied to the linear system with the Laplacian matrix. Each path attains a unique minimum and maximum which bound a region of inflection.

becomes closer to a diagonal matrix, and its condition number decreases, hence CG is more effective. Important to note is that our implementation does not entirely match this claim on the order of the CG method. Our `MMatrix` and `MVector` classes are not optimised for speed. For example, we are not optimised for multiplying with a sparse matrix, which would allow us to ignore all the entries containing 0.

## 2.4 Poisson's Equation in 2D

We are now interested in the problem in two variables  $\nabla^2 u(s, t) = 1$ , with boundary conditions  $u(s, 0) = u(s, 1) = u(0, t) = u(1, t) = 0$ . By the finite difference method, the matrix is of the form

$$(\mathbf{A})_{ij} = \begin{cases} 4 & \text{if } i = j \\ -1 & \text{if } |i - j| = n \\ -1 & \text{if } |i - j| = 1 \text{ and} \\ & (i + j) \bmod (2n) \neq 2n - 1 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

which we will refer to as the *Laplacian matrix*. Like before,  $\mathbf{b} = b_i = 1/(n + 1)^2$ . See Figure 3, which shows a selection of iterations for the 2D approximation to Poisson's equation. In Figures 4 and 5 we show the impacts of scaling the matrix size. Particularly in Figure 5, we have shown that our implementation is  $\mathcal{O}(d^3)$ , where  $d = n^2$  is the size of the system. Methods for solving a system of  $n$  linear equations are often  $\mathcal{O}(n^3)$  [2]. Solutions to Poisson's equation in 2D are shown in Figure

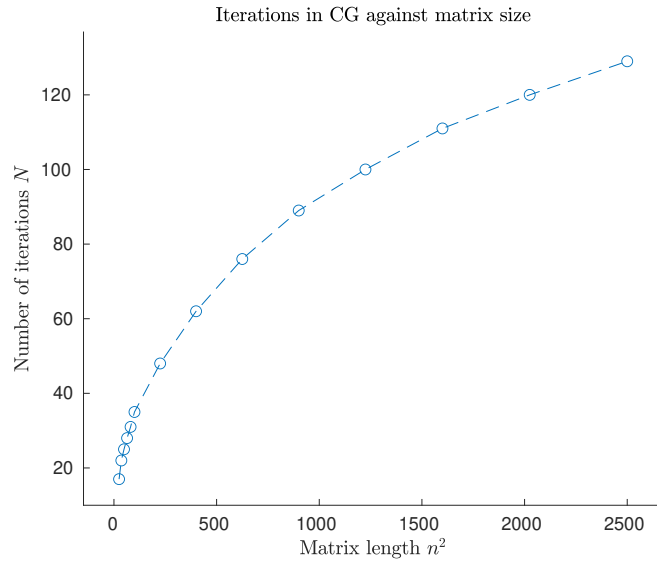


Figure 4: Number of iterations required for convergence for different matrix sizes. The curve appears logarithmic

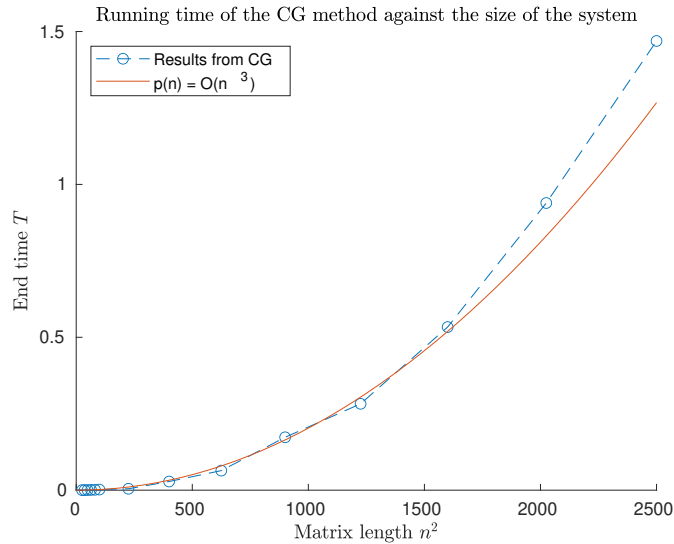


Figure 5: Running time, measured in seconds, of the C++ implementation of the CG method, against the size of the linear system we are solving. The polynomial approximation to our results was performed in MATLAB, and the fitting was unsuccessful for polynomials of degree  $< 3$ . The polynomial is technically  $\mathcal{O}(n^6)$  agreeing with our definition of  $n$  such that the system is  $n^2 \times n^2$ , but it is a degree 3 polynomial on the domain.

Computed solutions to Poisson's equation in 2D via the CG method

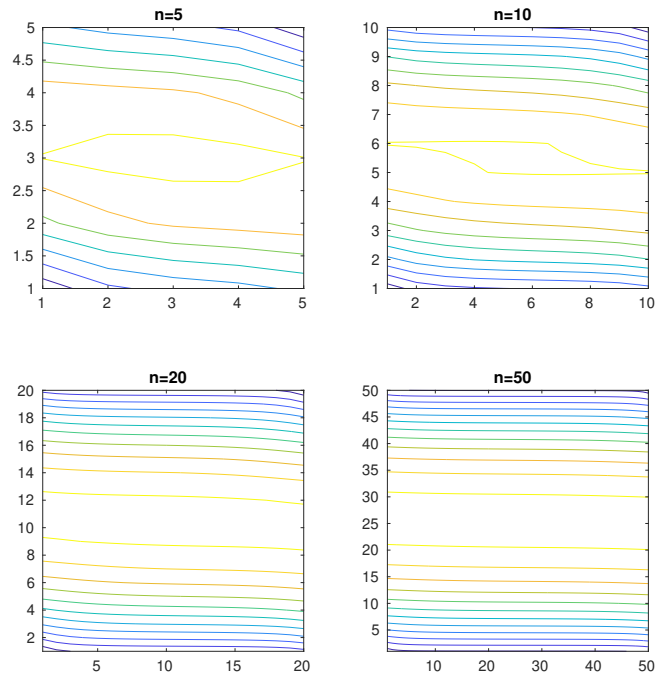


Figure 6: Solutions to Poisson's equation in 2D computed as a solution to a system of linear equations via the CG method. As  $n$  increases, our results converge and we see negligible change between  $n = 20$  and  $n = 50$ . Blue contours are lower in magnitude than the yellow curves near the centres of the images.

6, where we have shown a convergence of results as  $n$  increases. The boundary condition is attained on the horizontal boundaries.

## 2.5 Properties of the Conjugate Gradient Method

CG is only stable in exact arithmetic for symmetric, positive definite matrices [1]. As such, the modified inner product from earlier of a vector with itself  $\langle \mathbf{p}, \mathbf{p} \rangle_{\mathbf{A}} = \mathbf{p}^T \mathbf{A} \mathbf{p}$  is always positive. Furthermore if  $\mathbf{A}$  is not symmetric then

$$\begin{aligned} \langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{A}} &= \langle \mathbf{u}, \mathbf{A} \mathbf{v} \rangle \\ &= \langle \mathbf{A}^T \mathbf{u}, \mathbf{v} \rangle \\ &\neq \langle \mathbf{A} \mathbf{u}, \mathbf{v} \rangle \end{aligned}$$

is not an inner product.

All the stencil matrices we have looked at are symmetric (by inspection) and positive definite since they have entirely positive eigenvalues. This is a necessary criterion for CG to succeed. Floating point arithmetic may lead to rounding errors during computation, meaning that some methods will not succeed as they may in exact arithmetic.

Given the  $5 \times 5$  matrix

$$\mathbf{A} = \begin{bmatrix} 2 & -1 & 0 & 0 & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & -1 & 2 \end{bmatrix}$$

we can compute its inverse

$$\mathbf{A}^{-1} = \begin{bmatrix} \frac{5}{6} & \frac{2}{3} & \frac{1}{2} & \frac{1}{3} & \frac{1}{6} \\ \frac{2}{3} & \frac{4}{3} & \frac{1}{3} & \frac{2}{3} & \frac{1}{3} \\ \frac{1}{2} & \frac{1}{3} & \frac{3}{2} & \frac{1}{2} & \frac{1}{3} \\ \frac{1}{3} & \frac{2}{3} & \frac{1}{2} & \frac{4}{3} & \frac{2}{3} \\ \frac{1}{6} & \frac{1}{3} & \frac{1}{2} & \frac{2}{3} & \frac{5}{6} \end{bmatrix}$$

and find the condition number  $\|\mathbf{A}\|_F \|\mathbf{A}^{-1}\|_F \approx 20.7$ . The identity matrix has the minimum condition number 1. As we saw earlier, the condition number affects the time of the CG method. Particularly, Figure 2 showed that for larger diagonal entries, the method converges faster. For larger diagonal entries, the condition number is smaller. Hence a large condition number can be a strong factor affecting the run time of CG.

## 3 Conclusion

We can use the conjugate gradient method to solve a linear system  $\mathbf{A} \mathbf{x} = \mathbf{b}$  for a symmetric positive definite matrix. Uniquely, CG will always converge within the number of iterations that there are rows/columns in the matrix. It does this by computing search directions which are conjugate under  $\mathbf{A}$  and traversing them, computing how far each direction must be travelled by. Conjugate Gradient may be favourable to methods such as gradient descent or other options when we can guarantee stability. The time complexity of the method is improved for a sparse and also a well-conditioned matrix. Our results show the effectiveness of the CG method for solving finite difference approximations to differential equations in one and two dimensions, where we have been able to show accuracy and convergence of our results.

## References

- [1] Anne Greenbaum. “Chapter 4 - Effects of Finite Precision Arithmetic”. In: *Iterative methods for solving linear systems*. SIAM, 1997.
- [2] C. D. Meyer. “Linear Equations”. In: *Matrix analysis and applied linear algebra*. Society for Industrial and Applied Mathematics, 2000, pp. 14–16.
- [3] Jonathan Richard Shewchuk. “An introduction to the conjugate gradient method without the agonizing pain”. In: (1994), pp. 30–31.
- [4] Jonathan Richard Shewchuk. “An introduction to the conjugate gradient method without the agonizing pain”. In: (1994), pp. 37–38.

## A Appendix - C++ Code Implementations

### A.1 Compiled Code

#### A.1.1 main.cpp

```
1 #include <vector>
2 #include <algorithm>
3 #include <iostream>
4 #include <fstream>
5 #include <cmath>
6 #include <iomanip>
7 #include <string>
8 #include <chrono>
9 #include <cstdlib>
10 #include <ctime>
11 #include <ratio>
12
13 #include "mvector.h"
14 #include "mmatrix.h"
15
16 using namespace std;
17 using namespace chrono;
18
19 MVector conjugateGradientSolve(const MMatrix&, const MVector&, const MVector&, const
    string filename);
20 void computeLaplacian(int n, const string& filename);
21
22 int main()
23 {
24     MVector b(5, 1.0/36.0);
25
26     MMatrix A(5, 5);
27     makeTridiagonal(A, 2, -1);
28
29     // starting guess is elementary basis vector
30     MVector x_0 = {1,0,0,0,0};
31
32     //regular 5x5 solve (3.35-3.36)
33     MVector x = conjugateGradientSolve(A, b, x_0, "tridiagonal_convergence_0611.txt")
    ;
34     cout << x << endl;
35
36     MMatrix A1(5, 5), A2(5, 5), A3(5, 5);
```



```

37     makeTridiagonal(A1, 2+1, -1);
38     makeTridiagonal(A2, 2+4, -1);
39     makeTridiagonal(A3, 2+10, -1);
40
41     MVector b1(5, 2.5);
42
43     // perform CG for different values of alpha (3.37)
44     x = conjugateGradientSolve(A1, b1, x_0, "tridiagonal_a1_conv.txt");
45     x = conjugateGradientSolve(A2, b1, x_0, "tridiagonal_a4_conv.txt");
46     x = conjugateGradientSolve(A3, b1, x_0, "tridiagonal_a10_conv.txt");
47
48     MMatrix P(16, 16);
49     makeLaplacian(P, 4, -1);
50     cout << P << endl;
51
52     // Laplace's equation (3.39)
53     computeLaplacian(5, "lap_conv_n5_0611.txt");
54     computeLaplacian(6, "lap_conv_n6_0611.txt");
55     computeLaplacian(7, "lap_conv_n7_0611.txt");
56     computeLaplacian(8, "lap_conv_n8_0611.txt");
57     computeLaplacian(9, "lap_conv_n9_0611.txt");
58     computeLaplacian(10, "lap_conv_n10_0611.txt");
59     computeLaplacian(20, "lap_conv_n20_0611.txt");
60     computeLaplacian(50, "lap_conv_n50_0611.txt");
61
62
63     return 0;
64 }
65
66 // conjugate gradient method implementation
67 MVector conjugateGradientSolve(const MMatrix& A, const MVector& b, const MVector& x_0
    , const string filename = "iteration.txt")
68 {
69     // initialisations
70     int maxIterations = 1000;
71     double tolerance = 1e-6;
72
73     if (A.Rows() != A.Cols())
74     {
75         throw invalid_argument("error: matrix is not square");
76     }
77     int dim = A.Rows();
78
79     MVector rPrev(dim), pPrev(dim), xPrev(dim), r(dim), p(dim), x(dim);
80     double alpha, beta;
81
82     // solver writes a file, an optional filename argument is given
83     ofstream of;
84     int w = 12;
85     of.open(filename);
86     if (!of)
87     {
88         throw invalid_argument("error: failed to open file");
89     }
90     of << setw(w) << "Time t" << " | "
91     << setw(w) << "L2 norm" << " | "
92     << setw(w) << "LInf norm" << "\n";
93
94     xPrev = x_0;
95     rPrev = b - A*xPrev;

```

```

96     pPrev = rPrev;
97
98     high_resolution_clock::time_point t1 = high_resolution_clock::now();
99
100    of << setw(w) << 0 << " | "
101    << setw(w) << rPrev.L2Norm() << " | "
102    << setw(w) << rPrev.LInfNorm() << "\n";
103
104    for (int iter=0; iter<maxIterations; iter++)
105    {
106        // ...calculate new values for x and r here...
107        alpha = dot(rPrev, rPrev)/dot(pPrev, A*pPrev);
108        x = xPrev + alpha*pPrev;
109        r = rPrev - alpha*(A*pPrev);
110
111        // write the iteration
112        high_resolution_clock::time_point t2 = high_resolution_clock::now();
113        duration<double> time_span = duration_cast<duration<double>>(t2 - t1);
114
115        of << setw(10) << time_span.count() << " | "
116        << setw(10) << r.L2Norm() << " | "
117        << setw(10) << r.LInfNorm() << "\n";
118
119        // check if solution is accurate enough
120        if (r.L2Norm() < tolerance) break;
121
122        // ...calculate new conjugate vector p here...
123        beta = dot(r, r)/dot(rPrev, rPrev);
124        p = r + beta*pPrev;
125
126        // step the placeholder variables for the next iteration
127        pPrev = p;
128        rPrev = r;
129        xPrev = x;
130    }
131
132    of.close();
133    return x;
134 }
135
136 // procedure to compute a solution using the conjugate gradient method, given the
137 // value n and the file to write to
138 void computeLaplacian(int n, const string& filename)
139 {
140     MMatrix H(n*n, n*n);
141     makeLaplacian(H, 4, -1);
142
143     MVector a(n*n, 1.0/pow(n+1, 2));
144     // initial guess is always the unit vector
145     MVector x(n*n);
146     x[0] = 1;
147
148     conjugateGradientSolve(H, a, x, filename);
149 }

```

### A.1.2 mmatrix.cpp

```

1 #include <vector>
2 #include <iostream>
3 #include <iomanip>
4 #include <ios>

```

```

5 #include <algorithm>
6 #include <cmath>
7
8 #include "mmatrix.h"
9 #include "mvector.h"
10
11 using namespace std;
12
13 MVector operator*(const MMatrix& A, const MVector& v)
14 {
15     // check compatibility
16     if (A.Cols() != v.size())
17     {
18         cout << "wrong dims" << endl;
19         exception e;
20         throw e;
21     }
22     // initialise output vector
23     int vout_size = A.Rows();
24     MVector vout = MVector(vout_size);
25
26     for (int i = 0; i < A.Rows(); i++)
27     {
28         // sum across j
29         double sum = 0;
30         for (int j = 0; j < A.Cols(); j++)
31         {
32             sum += A(i, j)*v[j];
33         }
34         vout[i] = sum;
35     }
36     return vout;
37 }
38
39 ostream& operator<<(ostream& os, const MMatrix& A)
40 {
41     for (int i = 0; i < A.Rows(); i++)
42     {
43         os << "| ";
44         for (int j = 0; j < A.Cols(); j++)
45         {
46             os << setw(3) << A(i, j) << " ";
47         }
48         os << "\n";
49     }
50     return os;
51 }
52
53 void makeTridiagonal(MMatrix& A, double d, double bd)
54 {
55     /*
56     takes the arguments for the diagonal and band, and passes the matrix by reference.
57     A must already be the size we want, but its entries will be entirely overwritten.
58     */
59     if (A.Rows() != A.Cols())
60     {
61         throw invalid_argument("error: matrix not square");
62     }
63     for (int i = 0; i < A.Rows(); i++)
64     {

```

```

65     for (int j = 0; j < A.Cols(); j++)
66     {
67         if (i == j)
68         {
69             A(i, j) = d;
70         }
71         else if (abs(i - j) == 1)
72         {
73             A(i, j) = bd;
74         }
75         else
76         {
77             A(i, j) = 0;
78         }
79     }
80 }
81 }
82
83 void makeLaplacian(MMatrix& A, double d, double bd)
84 {
85     /*
86     behaves similarly to makeTridiagonal in its arguments. for solutions to the
87     Laplacian equation.
88     */
89     if (A.Rows() != A.Cols())
90     {
91         throw invalid_argument("error: matrix not square");
92     }
93     double sz = sqrt(A.Rows());
94     if (floor(sz) != sz)
95     {
96         // A must be n^2 by n^2
97         throw invalid_argument("error: matrix is not of dimension n^2.n^2");
98     }
99     int n = sz;
100    for (int i = 0; i < A.Rows(); i++)
101    {
102        for (int j = 0; j < A.Cols(); j++)
103        {
104            if (i == j)
105            {
106                A(i, j) = d;
107            }
108            else if (abs(i - j) == n)
109            {
110                A(i, j) = bd;
111            }
112            else if ((abs(i - j) == 1) && ((i + j) % 2*n != 2*n - 1))
113            {
114                A(i, j) = bd;
115            }
116            else
117            {
118                A(i, j) = 0;
119            }
120        }
121    }

```

### A.1.3 mvector.cpp

```
1 #include <vector>
2 #include <iostream>
3 #include <iomanip>
4 #include <algorithm>
5 #include <cmath>
6
7 #include "mvector.h"
8
9 using namespace std;
10
11 // constructors in header
12
13 MVector operator*(const double& lhs, const MVector& rhs)
14 {
15     MVector temp = rhs;
16     for (int i = 0; i < temp.size(); i++)
17     {
18         temp[i] *= lhs;
19     }
20     return temp;
21 }
22
23 MVector operator*(const MVector& lhs, const double& rhs)
24 {
25     return rhs*lhs;
26 }
27
28 MVector operator+(const MVector& lhs, const MVector& rhs)
29 {
30     if (lhs.size() != rhs.size())
31     {
32         exception e;
33         throw e;
34     }
35     MVector temp(lhs.size());
36     for (int i = 0; i < temp.size(); i++)
37     {
38         temp[i] = lhs[i] + rhs[i];
39     }
40     return temp;
41 }
42
43 MVector operator-(const MVector& lhs, const MVector& rhs)
44 {
45     return lhs+(-1*rhs);
46 }
47
48 // alternative overload such that given MVector v, -v is allowed
49 MVector MVector::operator-()
50 {
51     MVector vec(v.size());
52     for (int i = 0; i < vec.size(); i++)
53     {
54         vec[i] = -v[i];
55     }
56     return vec;
57 }
58
```

```

59 MVector operator/(const MVector& lhs, const double& rhs)
60 {
61     return (1/rhs)*lhs;
62 }
63
64 ostream& operator<<(ostream& os, const MVector& v)
65 {
66     int n = v.size();
67     os << "(";
68     for (int i = 0; i < n-1; i++)
69     {
70         os << setw(10) << v[i] << ", ";
71     }
72     os << v[n-1];
73     os << ")";
74     return os;
75 }
76
77 double MVector::LInfNorm() const
78 {
79     double maxAbs = 0;
80     size_t s = size();
81     for (int i=0; i<s; i++)
82     {
83         maxAbs = max(abs(v[i]), maxAbs);
84     }
85     return maxAbs;
86 }
87
88 double MVector::L2Norm() const
89 {
90     double sum = 0;
91     for (int i = 0; i < v.size(); i++)
92     {
93         sum += v[i]*v[i];
94     }
95     return sqrt(sum);
96 }
97
98 double dot(const MVector& lhs, const MVector& rhs)
99 {
100     if (lhs.size() != rhs.size())
101     {
102         throw invalid_argument("error: dot product not defined for vectors of non-equal
103             dimension");
104     }
105     double sum = 0;
106     for (int i = 0; i < lhs.size(); i++)
107     {
108         sum += lhs[i]*rhs[i];
109     }
110     return sum;
111 }

```

## A.2 Included Code

### A.2.1 mmatrix.h

```

1 #ifndef MMATRIX_H // the 'include guard'

```

```

2 #define MMATRIX_H
3
4 #include <vector>
5 #include <iostream>
6 #include "mvector.h"
7
8 using namespace std;
9
10 // Class that represents a mathematical matrix
11 class MMatrix
12 {
13 public:
14     // constructors
15     MMatrix() : nRows(0), nCols(0) {}
16     MMatrix(int n, int m, double x = 0) : nRows(n), nCols(m), A(n * m, x) {}
17
18     // set all matrix entries equal to a double
19     MMatrix &operator=(double x)
20     {
21         for (unsigned i = 0; i < nRows * nCols; i++) A[i] = x;
22         return *this;
23     }
24
25     // access element, indexed by (row, column) [rvalue]
26     double operator()(int i, int j) const
27     {
28         return A[j + i * nCols];
29     }
30
31     // access element, indexed by (row, column) [lvalue]
32     double &operator()(int i, int j)
33     {
34         return A[j + i * nCols];
35     }
36
37     friend MVector operator*(const MMatrix& A, const MVector& v);
38     friend ostream& operator<<(ostream& os, const MMatrix& A);
39
40     // size of matrix
41     int Rows() const { return nRows; }
42     int Cols() const { return nCols; }
43
44 private:
45     unsigned int nRows, nCols;
46     vector<double> A;
47 };
48
49 void makeTridiagonal(MMatrix& A, double d, double bd);
50 void makeLaplacian(MMatrix& A, double d, double bd);
51
52 #endif

```

## A.2.2 mvector.h

```

1 #ifndef MVECTOR_H // the 'include guard'
2 #define MVECTOR_H // see C++ Primer Sec. 2.9.2
3
4 #include <vector>
5

```

```

6 using namespace std;
7
8 // Class that represents a mathematical vector
9 class MVector
10 {
11 public:
12 // constructors
13 MVector() {}
14 MVector(int n) : v(n) {}
15 MVector(int n, double x) : v(n, x) {}
16 MVector(initializer_list<double> l) : v(l) {}
17
18 // access element (lvalue) (see example sheet 5, q5.6)
19 double &operator[](int index)
20 {
21     return v[index];
22 }
23
24 // access element (rvalue) (see example sheet 5, q5.7)
25 double operator[](int index) const {
26     return v[index];
27 }
28
29 // operator overloads
30 MVector operator-();
31
32 friend MVector operator*(const double& lhs, const MVector& rhs);
33 friend MVector operator*(const MVector& lhs, const double& rhs);
34 friend MVector operator+(const MVector& lhs, const MVector& rhs);
35 friend MVector operator-(const MVector& lhs, const MVector& rhs);
36 friend MVector operator/(const MVector& lhs, const double& rhs);
37
38 friend ostream& operator<<(ostream& os, const MVector& v);
39
40 // get size
41 int size() const { return v.size(); } // number of elements
42
43 //vector norms
44 double LInfNorm() const;
45 double L2Norm() const;
46
47
48 private:
49     vector<double> v;
50 };
51
52 // dot product
53 double dot(const MVector& lhs, const MVector& rhs);

```

## B Appendix - Extra Code

### B.1 Octave verification of linear solver

```

octave:11> A
A =
2 -1 0 0 0
-1 2 -1 0 0

```



```
0 -1 2 -1 0
0 0 -1 2 -1
0 0 0 -1 2
```

```
octave:12> b
```

```
b =
```

```
0.027778
0.027778
0.027778
0.027778
0.027778
```

```
octave:13> A\b
```

```
ans =
```

```
0.069444
0.111111
0.125000
0.111111
0.069444
```

## C Appendix - Notes on the problems addressed

The code provided does not immediately provide all of the results stated. I have enclosed the code as is due to time constraints. The current implementation writes a file which documents computation time for any execution of `conjugateGradientSolve()`. It is stated in the body that our code produces a file which consists of given results. The norm values are consistent, but the times will change due to different hardware and system resources if executed. The code was modified when it was needed to write the solutions to Poisson's equation in 2D in matrix form.

There is confusion with the problem for Poisson's equation in 2D. The problem states that we are solving  $\nabla^2 u = 1$  and the formulation for this is analogous to the 1D problem given, however for the 1D problem the model equation is  $u'(s) = -1$ . Our solutions for both are positive valued. If we are solving the 2D equation our solution should be negative, but the matrix and vector in the solution are verified to be correct as given by the problem statement. Furthermore, the solution seems to only satisfy the boundary conditions on the horizontal boundary lines.

All code compiled successfully on Ubuntu 22.04 Linux with g++ using `-O2` optimisation, faster floating point arithmetic `-ffast-math` and native hardware priority `-march=native`.