# Numerical Methods for CSE
# Final Exam - Summer 2019

**Prof. Rima Alaifari**

**ETH Zürich, D-MATH**

Exam Organisers:
Soumil Gurjar, Pratyuksh Bansal

August 26, 2019

**Exercise 1**. *Sparsity Preservation* (26 pts) [*Template: 1.cpp*]

Consider the following problem: Given $\mathbf{A} \in \mathbb{R}^{m,n}$, $m \geq n$, $\text{rank}(\mathbf{A}) = n$, $\mathbf{b} \in \mathbb{R}^m$, we want to obtain the unique least squares solution of the linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$.

(a) (2 pts) Implement a C++ function

```
VectorXd normalEquationSolve(const MatrixXd &A, const VectorXd &b);
```

that returns the least squares solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$ by solving the normal equations.

**Solution:**

```
17 VectorXd normalEquationSolve(const MatrixXd &A, const VectorXd &b){
18     if (b. size () != A.rows() ) throw std::runtime_error("Dimension mismatch");  //
       Cholesky solver
19     VectorXd x = (A.transpose()*A).llt().solve(A.transpose()*b);
20     return x ;
21 }
```

(b) (2 pts) Now consider a matrix $\mathbf{A}$ which is sparse. While performing the computation of the solution by normal equations, the sparsity is not necessarily preserved. Construct an example of a sparse $n \times n$ matrix such that the system matrix of the normal equations has no zero entries.

**Solution:** For example, considering a sparse matrix

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

when computing the solution by normal equations, we perform the matrix product $\mathbf{A}^\top\mathbf{A}$, which would be a dense matrix of only ones.

(c) (6 pts) Through an augmented system $\mathbf{C}\mathbf{y} = \mathbf{d}$, sparsity can be maintained. Implement two C++ functions,

```
SparseMatrix<double> computeC(SparseMatrix<double> &A, const VectorXd &b);
```

and

```
VectorXd computeRHS(SparseMatrix<double> &A, const VectorXd &b);
```

that respectively compute the *sparse* matrix $\mathbf{C}$ and RHS-vector $\mathbf{d}$.

**Solution:**

```
26 SparseMatrix<double> computeC(SparseMatrix<double> &A, const VectorXd &b){
27     int m = A.rows();
28     int n = A.cols();
29     int sizeC = m+n;
30     std :: vector < Triplet <double > > triplets_C ;
31     triplets_C.reserve(m+2*A.nonZeros()) ;
32
```

```
33      // Fill the triplets with elements from A and A_transpose
34      for(int i = 0; i < A.outerSize(); ++i){
35          for(SparseMatrix<double>::InnerIterator it(A,i); it; ++it){
36              triplets_C.emplace_back(it.row(),it.col()+m,it.value());
37              triplets_C.emplace_back(it.col()+m,it.row(),it.value());
38          }
39      }
40
41      //Fill the triplets corresponding to the Identity Block
42      for (int k=0; k<m; ++k){
43          triplets_C.push_back(Triplet<double>(k ,k ,-1.0));
44      }
45
46      SparseMatrix <double> C(sizeC, sizeC) ;
47      C.setFromTriplets(triplets_C.begin(), triplets_C.end());
48      C.makeCompressed();
49      return C;
50 }
```

```
53 VectorXd computeRHS(SparseMatrix<double> &A, const VectorXd &b){
54      int m = A.rows();
55      int n = A.cols();
56      int sizeC = m+n;
57      VectorXd d = VectorXd :: Zero(sizeC);
58      d.head(m) = b;
59
60      return d ;
61 }
```

(d) (4 pts) Implement a C++ function

```
VectorXd computeSolution(SparseMatrix<double> &A, const VectorXd &b);
```

that solves the modified system using a standard sparse solver (SparseLU) and returns the required solution vector $\mathbf{x}$ corresponding to the original system $\mathbf{Ax} = \mathbf{b}$.

**Solution:**

```
66 VectorXd computeSolution(SparseMatrix<double> &A, const VectorXd &b){
67      VectorXd x, y;
68      VectorXd d = computeRHS(A, b);
69      SparseMatrix<double> C = computeC(A, b);
70      // Perform sparse elimination
71      SparseLU<SparseMatrix<double>> solver(C);
72      y = solver.solve(d);
73
74      return x = y.tail(A.cols()); //The actual solution vector 'x' is the last 'n'
        elements of y
75 }
```

(e) (6 pts) Now we consider a specific example of the 1D Poisson equation, $-\ddot{u} = f$ on $\Omega = [0, 1]$, with Dirichlet boundary conditions $u(0) = u(1) = 0$. We apply a finite difference method, which allows us to obtain the solution to this problem by solving a linear system of equations. This is done in the following way:

The independent variable, $x \in \Omega$, is discretized by choosing an equidistant grid $\Omega_N = \{x_j\}_{j=0}^{N+1} \in \Omega$, where the grid points are given by $x_j = j\Delta x$, and where the spacing between the grid points is referred to as the mesh width, $\Delta x = \frac{1}{N+1}$. This means that $x_0 = 0$ is the left endpoint, and $x_{N+1} = 1$ is the right endpoint. In between, there are $N$ internal points in the interval $[0, 1]$.

The solution will be approximated numerically on this grid, by a vector $\mathbf{u} = \{u_j\}_1^N$ of internal points, augmented by the boundary values, $u_0 = 0$ and $u_{N+1} = 0$. The vector $\mathbf{u}$ approximates the solution to the differential equation, according to $u_j \approx u(x_j)$. Now, the second order derivative $-\ddot{u}$ is approximated by a symmetric finite difference scheme,

$$-\ddot{u}(x_j) = -\left(\frac{u_{j-1} - 2u_j + u_{j+1}}{(\Delta x)^2}\right).$$

Applying the finite difference method for all the interior points, obtain a linear system of equations of the form

$$L_N\mathbf{u} = \mathbf{f},$$

where the vector $\mathbf{f} = (f(x_1), f(x_2), \ldots, f(x_N))^\top$, $\mathbf{u} = (u_1, u_2, \ldots, u_N)^\top$ and $L_N \in \mathbb{R}^{N,N}$. Explicitly write out the matrix $L_N$.

**Solution:** A symmetric finite difference scheme leads to a linear system of equations:

$$\left(\frac{0 + 2u_1 - u_2}{(\Delta x)^2}\right) = f(x_1)$$

$$\left(\frac{-u_{j-1} + 2u_j - u_{j+1}}{(\Delta x)^2}\right) = f(x_j) \qquad j = 2 \ldots N - 1$$

$$\left(\frac{-u_{N-1} - 2u_N - 0}{(\Delta x)^2}\right) = f(x_N),$$

where the first and last equations include the zero-valued boundary conditions, and the middle equation is the generic equation, only involving internal points.

This linear system is conveniently represented in matrix-vector form $L_N\mathbf{u} = \mathbf{f}$, as

$$\frac{1}{(\Delta x)^2}
\begin{pmatrix}
2 & -1 & 0 & 0 & \cdots & & \cdots & 0 \\
-1 & 2 & -1 & 0 & \cdots & & & 0 \\
0 & -1 & 2 & -1 & \ddots & & & \\
\vdots & & \ddots & \ddots & \ddots & & & \\
& & & \ddots & \ddots & & \ddots & \\
& & & & -1 & 2 & -1 & 0 \\
\vdots & & & & & -1 & 2 & -1 \\
0 & \cdots & & & & 0 & -1 & 2
\end{pmatrix}
\begin{pmatrix}
u_1 \\ \vdots \\ \\ \\ \\ \vdots \\ u_N
\end{pmatrix}
=
\begin{pmatrix}
f(x_1) \\ \vdots \\ \\ \\ \\ \vdots \\ f(x_N)
\end{pmatrix}$$

Thus, the $N \times N$ tridiagonal matrix can be written

$$L_N = \frac{1}{(\Delta x)^2}\mathbf{T}$$

with

$$\mathbf{T} = tridiag(-1 \quad 2 \quad -1).$$

(f) (6 pts) For the 1D Poisson problem described above, if $f = \pi^2 \sin(\pi x)$ and $N = 25, 50, 100, ..., 1600$, implement two functions, `constructDenseA_Poisson` and `constructSparseA_Poisson`, that construct the matrix $L_N$ in dense and sparse format respectively for each mesh-width. Solve the corresponding linear system $L_N \mathbf{u} = \mathbf{f}$, using the normal equations from sub-problem (a), as well as the modified system from sub-problem (d). Compare the two solutions to verify your implementation, and also compare their run-times for each value of $N$. Explain your observations.

Hint: A dense matrix `A_dense` can be converted into a sparse matrix `A_sparse` by using `A_sparse = A_dense.sparseView()`. Moreover, the code to compare the solutions and run-times is already implemented within the template.

**Solution:**

The exact solution to this problem is $u = \sin(\pi x)$. Thus, at the interior nodes, we have the exact solution $u_i = \sin(\pi x_i)$

```
83  //Construct dense matrix 'A' obtained through Finite Difference Method for the
        poisson problem
84  MatrixXd constructDenseA_Poisson(int N){
85      // Mesh width
86      double delta_x = 1.0/(N+1);
87      MatrixXd A_dense = MatrixXd::Zero(N,N);
88
89      A_dense.diagonal() = VectorXd::Constant(N,2);
90      A_dense.diagonal<1>() = VectorXd::Constant(N-1,-1);
91      A_dense.diagonal<-1>() = VectorXd::Constant(N-1,-1);
92      A_dense *= 1/(delta_x*delta_x);
93
94      return A_dense;
95  }
96
97  //Construct sparse matrix 'A' obtained through Finite Difference Method for the
        poisson problem
98  SparseMatrix <double> constructSparseA_Poisson(MatrixXd &A_dense, int N){
99      SparseMatrix <double> A_sparse = A_dense.sparseView(); //Convert dense matrix A
        into a sparse form
100     return A_sparse;
101 }
```

<div align="center">question_1/codes/1_sol.cpp</div>

RunTime Comparison: For $N \geq 75$, we observe that the run-time for the augmented system method is lower than that of the normal system method. The benefit of the augmented system and preservation of sparsity can be clearly seen for the cases with $N = 200, 400, 800, 1600$.

| N | maxerr_normal | runTime_normal | maxerr_sparse | runTime_sparse |
|---|---|---|---|---|
| 25 | 0.00121756 | 0.000444536 | 0.00121756 | 0.00176432 |
| 50 | 0.000316122 | 0.00176937 | 0.000316122 | 0.00267536 |
| 75 | 0.000142406 | 0.00444235 | 0.000142406 | 0.00380388 |
| 100 | 8.06202e-05 | 0.00890637 | 8.06203e-05 | 0.00535135 |
| 200 | 2.03575e-05 | 0.0652939 | 2.03572e-05 | 0.00960842 |
| 400 | 5.12185e-06 | 0.467776 | 5.11479e-06 | 0.0188642 |
| 800 | 1.20004e-06 | 3.66319 | 1.2819e-06 | 0.0398366 |
| 1600 | 4.1513e-05 | 29.0481 | 3.20874e-07 | 0.0760622 |

Figure 1: Run-time comparison between Normal Equation method and Augmented System method

**Exercise 2.** *Interpolation* (28 pts) [*Template: 2.cpp*]

Consider the function

$$f(x) = \begin{cases} f_1(x) := 2 + \frac{1}{1+25(6x-1)^2}, & x \in [0, \frac{1}{3}), \\ f_2(x) := 1 + \cos(\pi x), & x \in [\frac{1}{3}, 1]. \end{cases} \tag{1}$$

Let $I_{\mathcal{T}_N}^p[f]$ be a piecewise Lagrange interpolant of the function $f$, for a mesh $\mathcal{T}_N$ on the interval $[0, 1]$. The mesh $\mathcal{T}_N$ has $N$ cells and the local polynomial degree $p \in \mathbb{N}_0$ of the interpolant is the same for each cell in the mesh.

Two functions are provided, see `2.cpp`, to generate equidistant and Chebyshev interpolation nodes in the reference interval $[0, 1]$. Let $t = \{t_0, t_1, \ldots, t_p\} \subset [0, 1]$ denote the set of these reference nodes.

(a) (2 pts) Is $f \in C^0([0, 1])$? Explain.

**Solution:** No, $f \notin C^0([0, 1])$. The given function is discontinuous at $x = 1/3$,

$$f_1(\frac{1}{3}) = 2 + \frac{1}{26} \approx 2.04 \neq 1.5 = f_2(\frac{1}{3}).$$

It is continuous for $x \in [0, 1] \setminus \{\frac{1}{3}\}$.

(b) (5 pts) For $N = 2$, design the interpolant:

  i) specify the mesh and

  ii) specify the strategy to generate local interpolation nodes,

such that the interpolation error $\|I_{\mathcal{T}_N}^p[f] - f\|_{L^\infty([0,1])}$ is small for an arbitrarily large $p$. Support your answer.

**Solution:**

  i) For $N = 2$, we choose the mesh vertices $\{0, \frac{1}{3}, 1\}$. For any other choice of mesh vertices $\{0, \xi, 1\}, \xi \neq \frac{1}{3}$, the point of discontinuity $x = \frac{1}{3}$ is included in either $[0, \xi]$ or $[\xi, 1]$. Therefore, the function is not smooth in one of the two cells and increasing the polynomial degree $p$ of the interpolation will not reduce the interpolation error.

  ii) In the cell $(\frac{1}{3}, 1)$, we can simply generate *equidistant* local interpolation nodes. In the cell $(0, \frac{1}{3})$, the given function is an affine transformation of the **Runge's function**, thus, Lagrange interpolation is unstable for equidistant nodes but stable for *Chebyshev* nodes, for a large polynomial degree $p$.

(c) (2 pts) Consider a *uniform* mesh $\mathcal{T}_M$ on the interval $[a, b]$. Implement a function

```
Eigen::MatrixXd genNodes(const double a,
                         const double b,
                         const int M,
                         const Eigen::VectorXd &t);
```

that uses a vector of reference nodes $t \in \mathbb{R}^n$ and generates a matrix $\mathbb{R}^{n \times M}$ of interpolation nodes in the mesh $\mathcal{T}_M$.

*Hint:* Map the reference nodes to a given cell.

**Solution:**

```
133  Eigen::MatrixXd genNodes(const double a, const double b, const int M, const Eigen::
         VectorXd &knots)
134  {
135      const int n = knots.size();
136      double H = (b-a)/M;
137
138      Eigen::MatrixXd nodes(n,M);
139      for (int j=0; j<M; j++) {
140          nodes.col(j) = (a + j*H) + knots.array()*H;
141      }
142
143      return nodes;
144  }
```

question_2/codes/2_sol.cpp

(d) (8 pts) Given the declaration of a **class Newton**, implement the member functions:

   i) `Interpolate` : Computes the coefficients of a piecewise Lagrange interpolant in the Newton basis.

   ii) `Evaluate` : Uses the *Horner scheme* to evaluate a piecewise Lagrange interpolant in the Newton basis.

**Solution:**

```
154  void Newton::Interpolate(const Eigen::MatrixXd &y)
155  {
156      int n = nodes.rows();
157      int M = nodes.cols();
158      assert(y.rows() == n);
159      assert(y.cols() == M);
160
161      // set-up and solve linear system
162      for (int l=0; l<M; l++)
163      {
164          Eigen::MatrixXd A(n,n);
165          A.setZero();
166
167          for (int i=0; i<n; i++) {
168              A(i,0) = 1;
169              for (int j=1; j<=i; j++) {
170                  A(i,j) = A(i,j-1)*(nodes(i,l) - nodes(j-1,l));
171              }
172          }
173
174          alpha.col(l) = A.triangularView<Eigen::Lower>().solve(y.col(l));
175      }
176
177      // solve with divided difference method
178      /*for (int l=0; l<M; l++) {
179
180          alpha.col(l) = y.col(l);
181          for (int j=0; j<n-1; ++j) {
182              for (int i=n-1; i>j; --i) {
183                  alpha(i,l) = (alpha(i,l) - alpha(i-1,l)) / (nodes(i,l) - nodes(i-1-j,
         l));
```

```
184                }
185            }
186        }*/
187 }
```
question_2/codes/2_sol.cpp

```
196 double Newton::Evaluate(const int l, double x)
197 {
198     int n = alpha.rows();
199     double y = alpha(n-1,l);
200     for (int i = n-2; i >= 0; --i) { // Horner scheme
201         y = y * (x - nodes(i,l)) + alpha(i,l);
202     }
203
204     return y;
205 }
```
question_2/codes/2_sol.cpp

(e) (8 pts) Use the components from sub-problems (c) and (d), and write a C++ function to compute a piecewise Lagrange interpolant of the given function $f(x)$ from (1) on a uniform mesh $\mathcal{T}_N$. This function should also evaluate and return the interpolation error.

*Hint:* Use the given member function EvalError, of class **Newton**, to compute the interpolation error.

**Solution:**

```
209 // Given function f
210 auto f = [](const double x)
211            {
212                if (x < 1./3.)
213                    return (1.0 / (1.0 + 25.0 * (6*x-1) * (6*x-1)));
214                else
215                    return (1.0 + cos(M_PI*x));
216            };
217
218 // function f1
219 auto f1m = [](const Eigen::MatrixXd x)
220              {
221                  return (1.0 / (1.0 + 25.0 * (6*x.array()-1) * (6*x.array()-1))).
222     matrix();
223              };
224
224 // function f2
225 auto f2m = [](const Eigen::MatrixXd x)
226              {
227                  return (1.0 + cos(M_PI*x.array())).matrix();
228              };
229
230 ////TASK E
231 // Computes the piecewise Lagrange interpolant
232 // and evaluates the interpolation error
233 double piecewiseLagrange_cheb(const int p, const int l)
234 {
235     double a = 0;
236     double c = 1./3.;
```

```cpp
237        double b = 1;
238
239        int N1 = std::pow(2,l);
240        int N2 = 2*N1;
241
242        // generate grid
243        Eigen::VectorXd grid = Eigen::VectorXd::LinSpaced(N1+N2+1, a, b);
244
245        // generate reference nodes
246        Eigen::VectorXd knotsCheb = genKnotsCheb(p+1);
247        Eigen::VectorXd knotsEqui = genKnotsEqui(p+1);
248
249        // generate nodes
250        Eigen::MatrixXd x(p+1, N1+N2);
251        x.block(0, 0, p+1, N1) = genNodes(a, c, N1, knotsCheb);
252        x.block(0, N1, p+1, N2) = genNodes(c, b, N2, knotsEqui);
253
254        // generate data
255        Eigen::MatrixXd y(p+1, N1+N2);
256        y.block(0, 0, p+1, N1) = f1m(x.block(0, 0, p+1, N1));
257        y.block(0, N1, p+1, N2) = f2m(x.block(0, N1, p+1, N2));
258
259        // solve interpolant
260        Newton ipolf(grid, x);
261        ipolf.Interpolate(y);
262
263        // compute interpolation error
264        double error = ipolf.EvalError(f);
265
266        return error;
267    }
```

<div align="center">question_2/codes/2_sol.cpp</div>

or if equidistant nodes are used:

```cpp
270    double piecewiseLagrange_equi(const int p, const int l)
271    {
272        double a = 0;
273        double c = 1./3.;
274        double b = 1;
275
276        int N1 = std::pow(2,l);
277        int N2 = 2*N1;
278
279        // generate grid
280        Eigen::VectorXd grid = Eigen::VectorXd::LinSpaced(N1+N2+1, a, b);
281
282        // generate reference nodes
283        Eigen::VectorXd knotsEqui = genKnotsEqui(p+1);
284
285        // generate nodes
286        Eigen::MatrixXd x(p+1, N1+N2);
287        x.block(0, 0, p+1, N1) = genNodes(a, c, N1, knotsEqui);
288        x.block(0, N1, p+1, N2) = genNodes(c, b, N2, knotsEqui);
289
290        // generate data
```

```
291    Eigen::MatrixXd y(p+1, N1+N2);
292    y.block(0, 0, p+1, N1) = f1m(x.block(0, 0, p+1, N1));
293    y.block(0, N1, p+1, N2) = f2m(x.block(0, N1, p+1, N2));
294
295    // solve interpolant
296    Newton ipolf(grid, x);
297    ipolf.Interpolate(y);
298
299    // compute interpolation error
300    double error = ipolf.EvalError(f);
301
302    return error;
303 }
```

question_2/codes/2_sol.cpp

(f) (3 pts) Implement a `C++` function to study the convergence of your interpolant $I^p_{\mathcal{T}_N}$ with respect to $p$, for $p = 1, 2, 3, \ldots, 20$. Explain the results from your convergence test.

**Solution:**

```
308 void convergenceWrtp()
309 {
310     int l=0;
311
312     std::cout << "\n\nEquidistant nodes, convergence with respect to p:" << std::endl
        ;
313     for (int p=1; p<21; p++)
314     {
315         double error = piecewiseLagrange_equi(p, l);
316         std::cout << p << "\t" << error << std::endl;
317     }
318
319     std::cout << "\n\nChebyshev nodes, convergence with respect to p:" << std::endl;
320     for (int p=1; p<21; p++)
321     {
322         double error = piecewiseLagrange_cheb(p, l);
323         std::cout << p << "\t" << error << std::endl;
324     }
325 }
```

question_2/codes/2_sol.cpp

If equidistant nodes are used to interpolate the given function in the range $[0, 1/3]$, we do not observe convergence with respect to the local polynomial degree $p$ because of the **Runge phoenomenon**. Whereas, when the Chebyshev points are used, the interpolation scheme converges with respect to $p$.

**Exercise 3**. *Singly Diagonally Implicit Runge-Kutta Methods* (26 pts) [*Template: 3.cpp*]

For general implicit Runge-Kutta methods, the $k_i$'s cannot be evaluated successively since they are coupled in the system of implicit equations that is given for their determination. One way to decouple a non-linear system is to use diagonally implicit Runge-Kutta methods. A special family of those methods are Singly Diagonally Implicit Runge-Kutta methods (SDIRK), in which the matrix of the method is lower triangular and all the diagonal entries are equal. Continuing in this direction, we define the following one parameter family of SDIRK-methods:

$$
\begin{array}{c|cc}
\gamma & \gamma & 0 \\
1-\gamma & 1-2\gamma & \gamma \\
\hline
& \frac{1}{2} & \frac{1}{2}
\end{array}
$$

Table I: Singly Diagonally Implicit Runge-Kutta

(a) (3 pts) Determine the value(s) of $\gamma$ for which the corresponding Runge-Kutta method is of at least (consistency) order 3.

Hint: In addition to conditions necessary for a consistency order 2, the following conditions are necessary for a consistency order 3:

$$
\sum_{i=1}^{s} b_i c_i^2 = \frac{1}{3} \qquad \text{and} \qquad \sum_{i=1}^{s}\sum_{j=1}^{s} b_i a_{ij} c_j = \frac{1}{6},
$$

where $\{a_{ij}\}_{i,j=1}^{s}$, $\{b_i\}_{i=1}^{s}$, $\{c_j\}_{j=1}^{s}$ are entries of a standard Butcher tableau.

**Solution:** The following equations have to be satisfied to obtain a consistency order of at least 3:

$$
\sum_{i=1}^{s} b_i = 1
$$

$$
\sum_{i=1}^{s} b_i c_i = \frac{1}{2}
$$

$$
\sum_{i=1}^{s} b_i c_i^2 = \frac{1}{3}
$$

$$
\sum_{i=1}^{s}\sum_{j=1}^{s} b_i a_{ij} c_j = \frac{1}{6}
$$

Therefore, we have:

$$
\frac{1}{2} + \frac{1}{2} = 1
$$

$$
\frac{1}{2}\gamma + \frac{1}{2}(1-\gamma) = \frac{1}{2}
$$

$$
\frac{1}{2}\gamma^2 + \frac{1}{2}(1-\gamma)^2 = \frac{1}{3}
$$

$$
\frac{1}{2}\gamma^2 + 0 + \frac{1}{2}\gamma(1-2\gamma) + \frac{1}{2}\gamma(1-\gamma) = \frac{1}{6}
$$

Therefore, we have that the method will be of order 3 provided $\gamma$ satisfies the equation

$$
\gamma - \gamma^2 = \frac{1}{6}
$$

that is, for,

$$\gamma = \frac{3 \pm \sqrt{3}}{6}$$

(b) (4 pts) Write down the iterations of the Runge-Kutta method defined in Table I, when applied to the ODE

$$\dot{y} = \lambda y, \qquad y(0) = 1, \qquad \lambda \in \mathbb{C}$$

HINT: Use the stability function and the fact that the ODE is linear.

**Solution:** Notice that we have

$$y_{n+1} = S(\lambda h)y_n$$

where $S(z)$ is the stability function of the Runge-Kutta method.

The stability function of any Runge-Kutta method is given by

$$S(z) := \frac{\det\left(\mathbf{I} - z\mathbf{A} + z\mathbf{1}\mathbf{b}^\top\right)}{\det\left(\mathbf{I} - z\mathbf{A}\right)}.$$

Therefore, setting $\mathbf{A}$ and $\mathbf{b}$ as given in table I, we have

$$\det\left(\mathbf{I} - z\mathbf{A}\right) = (1 - z\gamma)^2$$

and

$$\det\left(\mathbf{I} - z\mathbf{A} + z\mathbf{1}\mathbf{b}^\top\right) = 1 + z(1 - 2\gamma) + z^2\left(\frac{1}{2} - 2\gamma + \gamma^2\right)$$

So

$$y_{n+1} = \frac{1 + \lambda h(1 - 2\gamma) + (\lambda h)^2\left(\frac{1}{2} - 2\gamma + \gamma^2\right)}{(1 - \lambda h\gamma)^2}y_n$$

(c) (3 pts) For which value(s) of $\gamma$ obtained from sub-problem (a) (the value(s) with which the method is of order 3), can we conjecture the corresponding Runge-Kutta method to be A-stable? Explain.

Hint: You may use the plotting script *'stabilityRegion_plot.cpp'* to visualize the region of stability for a chosen value of $\gamma$.

**Solution:** From fig. 2 and fig. 3, we can conjecture that the method is A-stable only for $\gamma = \frac{3+\sqrt{3}}{6}$ as it is stable in the entire left half of the complex plane ($|S(z)| < 1 \ \forall z \in \mathbb{C}^-$).

```
10 ////TASK C
11 // This function computes |S(z)| for a chosen value of 'gamma'
12 double computeSMag(std::complex<double> z, double gamma){
13     double Smag;
14     std::complex<double> Stab_func;
15
16     //To-do [task (c)]: Write stability function in terms of 'z' and gamma
17     Stab_func = (1.0+z*(1.0-2.0*gamma)+z*z*(0.5-2.0*gamma+pow(gamma,2)))/(pow((1.0-z*
        gamma),2));
18
19     Smag = std::abs(Stab_func);
20
21     return Smag;
22 }
23 ////END OF TASK C
24
```
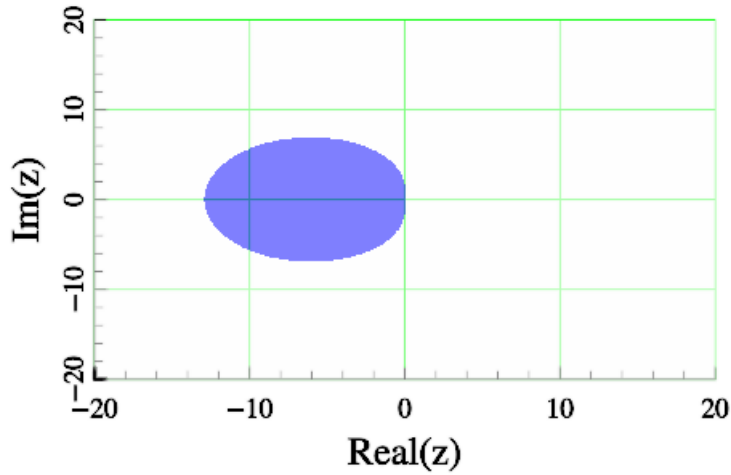
# Blue Region: |S(z)| < 1



Figure 2: Stability region for $\gamma = \frac{3-\sqrt{3}}{6}$ (Blue is stable region)

```cpp
int main() {
    //To-do [task (c)]: : Modify value for gamma in next line
    double gamma = (3.0-sqrt(3.0))/6.0;
    //double gamma = (3.0+sqrt(3.0))/6.0;
    //double gamma = 1.0;

    const std::complex<double> I(0.0, 1.0); //complex number 'iota'=sqrt(-1)
    double x_min = -20.0;
    double x_max = 20.0;
    double y_min = -20.0;
    double y_max = 20.0;
    int N_x = 1001;
    int N_y = 1001;
    double delta_x = (x_max - x_min)/(N_x-1);
    double delta_y = (y_max - y_min)/(N_y-1);
    VectorXd x_grid = VectorXd::LinSpaced(N_x, x_min, x_max);
    VectorXd y_grid = VectorXd::LinSpaced(N_y, y_min, y_max);

    // generate stability curve
    MatrixXcd z_grid(N_x,N_y);
    MatrixXd Smag_grid(N_x,N_y);
    MatrixXd Smag_grid_boolean(N_x,N_y);
    Smag_grid_boolean.setOnes();
    for(int k=0; k < N_y; ++k){
        for(int j=0; j < N_x; ++j){
            z_grid(j,k) = x_grid(j)+I*y_grid(k); // z = x + iota * y;
            Smag_grid(j,k) = computeSMag(z_grid(j,k),gamma); // Compute |S(z)| for
    given 'gamma'
            if (Smag_grid(j,k) >= 1.0) Smag_grid_boolean(j,k) = 0.0; //Create a
    boolean variable that takes value = 0.0 wherever |S(z)| >= 1 and takes value =
    1.0 wherever |S(z)| < 1.
        }
    }
```
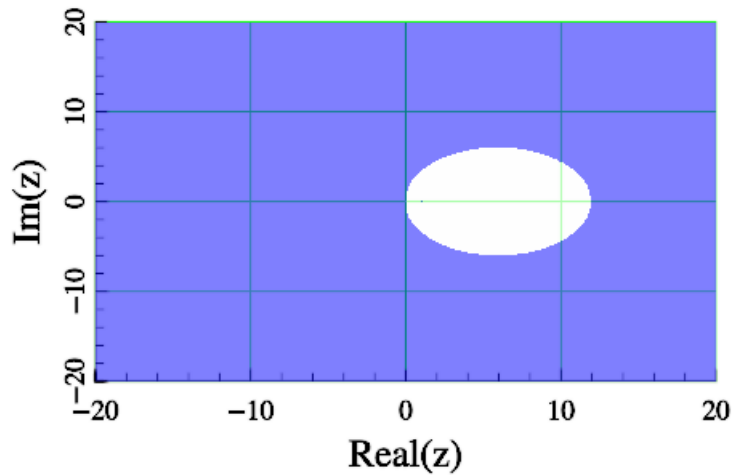
Figure 3: Stability region for $\gamma = \frac{3+\sqrt{3}}{6}$ (Blue is stable region)

```
56    // Plot stability function using MathGL (Blue region indicates region of |S(z)| <
      1 )
57    mglData x_dat(x_grid.data(), x_grid.size());
58    mglData y_dat(y_grid.data(), y_grid.size());
59    mglData dat(Smag_grid_boolean.rows(), Smag_grid_boolean.cols(), Smag_grid_boolean
      .data());
60    mglGraph gr;
61    gr.Title("Blue Region: |S(z)| < 1");
62    gr.SetRanges(x_min,x_max,y_min,y_max);
63    gr.Grid("xy","g");
64    gr.Axis(); gr.Label('x',"Real(z)",0); gr.Label('y',"Im(z)",0);
65    gr.Alpha(true);
66    gr.SetRange('c',0.0,1.0);
67    gr.Tile(x_dat,y_dat,dat,"wb");
68    gr.WriteFrame("stabilityPlot1.png");
69    //gr.WriteFrame("stabilityPlot2.png");
70
71    return 0;
```

question_3/codes/stabilityRegion_plot_sol.cpp

(d) (2 pts) We consider a scalar linear initial value problem of second order

$$\ddot{y} + \dot{y} + y = 0, \qquad y(0) = 1, \quad \dot{y}(0) = 0 \tag{2}$$

that should be solved numerically using the SDIRK-method described in Table I.

Formulate (2) as an initial value problem (IVP) for a linear first order system.

**Solution:**   Define $z_1 = y$, $z_2 = \dot{y}$, then the initial value problem

$$\ddot{y} + \dot{y} + y = 0, \qquad y(0) = 1, \quad \dot{y}(0) = 0 \tag{3}$$

is equivalent to the first order system

$$\dot{z}_1 = z_2$$
$$\dot{z}_2 = -z_1 - z_2, \tag{4}$$

with initial values $z_1(0) = 1$, $z_2(0) = 0$.

(e) (6 pts) Implement a C++/Eigen function

```
template <class StateType>
StateType sdirkStep(const StateType & z0, double h, double gamma);
```

that realizes the numerical evolution of one step of the SDIRK-method for the IVP obtained in sub-problem (d), starting from the initial value z0 and returning the value after a time step of size $h$.

Hint: For an SDIRK-method, compute stage $k_1$ first and then use it to compute stage $k_2$.

**Solution:** Let

$$\mathbf{A} := \begin{pmatrix} 0 & 1 \\ -1 & -1 \end{pmatrix}.$$

Then

$$k_1 = \mathbf{A}\mathbf{y}_0 + h\gamma\mathbf{A}k_1$$
$$k_2 = \mathbf{A}\mathbf{y}_0 + h(1 - 2\gamma)\mathbf{A}k_1 + h\mathbf{A}k_2,$$

so

$$k_1 = (\mathbf{I} - h\gamma\mathbf{A})^{-1}\mathbf{A}\mathbf{y}_0$$
$$k_2 = (\mathbf{I} - h\gamma\mathbf{A})^{-1}(\mathbf{A}\mathbf{y}_0 + h(1 - 2\gamma)\mathbf{A}k_1).$$

```
template <class StateType>
StateType sdirkStep(const StateType & z0, double h, double gamma) {
    // Matrix A for evaluation of f
    Eigen::Matrix2d A;
    A << 0.,1.,-1.,-1.;

    // Reuse factorization
    auto A_lu = (Eigen::Matrix2d::Identity()-h*gamma*A).partialPivLu();
    Eigen::Vector2d az = A*z0;

    // Increments
    Eigen::Vector2d k1 = A_lu.solve(az);
    Eigen::Vector2d k2 = A_lu.solve(az + h*(1-2*gamma)*A*k1);

    // Next step
    StateType z1 = z0 + h*0.5*(k1 + k2);
    return z1;
}
```

question_3/codes/3_sol.cpp

(f) (8 pts) Conduct a numerical experiment to deduce the convergence order of the method by computing the global error of the *A-stable* SDIRK method for the first order IVP from sub-problem (d) at the end time. Choose T=10 as end time and N=20,40,80,...,10240 as steps. Report the observed order of convergence.

**Solution:**

```
44  ////TASK F
45  template <class StateType>
46  std::vector<StateType> sdirkSolve(const StateType & z0, unsigned int N, double T,
        double gamma) {
47      // Solution vector
48      std::vector<StateType> res(N+1);
49      // Equidistant step size
50      const double h = T/N;
51
52      // Push initial data
53      res.at(0) = z0;
54
55      // Main loop
56      for(unsigned int i = 1; i <= N; ++i) {
57          res.at(i) = sdirkStep(res.at(i-1), h, gamma);
58      }
59
60      return res;
61  }
62  ////END OF TASK F
63
64  int main() {
65      // Initial data z0 = [y(0), y'(0)]
66      Eigen::Vector2d z0;
67      z0 << 1,0;
68      // Final time
69      const double T = 10;
70      // TO-DO task (f): Modify value of parameter gamma such that method is A-stable
71      const double gamma = (3.+std::sqrt(3.)) / 6.;
72      // Mesh sizes
73      std::vector<int> N = {20, 40, 80, 160, 320, 640, 1280, 2560, 5120, 10240};
74      // Exact solution (only y(t)) given initial data  z0 and t
75      auto yex = [&z0] (double t) {
76          return 1./3.*std::exp(-t/2.) * ( 3.*z0(0) * std::cos( std::sqrt(3.)*t/2. ) +
77          std::sqrt(3.)*z0(0) * std::sin( std::sqrt(3.)*t/2. ) +
78          2.*std::sqrt(3.)*z0(1) * std::sin( std::sqrt(3.)*t/2. ) );
79      };
80
81      // Store old error for rate computation
82      double errold = 0;
83      std::cout << std::setw(15) << "n" << std::setw(15) << "maxerr" << std::setw(15)
        << "rate" << std::endl;
84      // Loop over all meshes
85      for(unsigned int i = 0; i < N.size(); ++i) {
86          int n = N.at(i);
87          // Get solution
88          auto sol =  sdirkSolve(z0, n, T, gamma);
89          // Compute error
90          double err = std::abs(sol.back()(0) - yex(T));
91
92          // I/O
93          std::cout << std::setw(15) << n << std::setw(15) << err;
94          if(i > 0) std::cout << std::setw(15) << std::log2(errold /err);
95          std::cout << std::endl;
96
```

```
97          // Store old error
98          errold = err;
99      }
100 }
```

question_3/codes/3_sol.cpp

As seen in fig. 4, the numerically estimated order of convergence is 3, which is consistent with the expected order from sub-problem (a).

| n | maxerr | rate |
|---:|---|---|
| 20 | 0.000283567 | |
| 40 | 5.05272e-05 | 2.48856 |
| 80 | 7.73715e-06 | 2.70719 |
| 160 | 1.07013e-06 | 2.85402 |
| 320 | 1.40565e-07 | 2.92848 |
| 640 | 1.80048e-08 | 2.96478 |
| 1280 | 2.27798e-09 | 2.98255 |
| 2560 | 2.86467e-10 | 2.99132 |
| 5120 | 3.59161e-11 | 2.99567 |
| 10240 | 4.49619e-12 | 2.99785 |

Figure 4: Observed rate of convergence for SDIRK method