

SORRY PAUL

ARDEN RASMUSSEN

| |
|-------------------|
| Error |
| Remove |
| Appendix/Citation |
| Addition |
| Revise |

ABSTRACT. Finite Element Methods is a process that is utilized to numerically approximate solutions to other wise impossible partial differential equations provided with arbitrairy boundaries. This paper serves as a brief explination and introduction to the methods of finite element analysis. This paper provides a complete explanation of the every process involved, a reader should be able to implement all the necessary steps required in the analysis.

1. INTRODUCTION

The Method of Finite Elements (FEM) is a process that is utilized for numerically computing and apprximate solution to partial differential equations, on some specified domain. There are only a few cases where this process is analytically solvable, but for most situations this approximation is as close as technology currently allows us to achieve.

The uses of finite element methods in computation and analysis is often called Finite Element Anaylysis (FEA), while the mathmatics of the process is called FEM.

We will utilize the process closly set forth by [KH15]. This method is constructed around the Galerkin method for generating the approximation. The scope of this introduction is limited to two dimensions, but all of the mathematics is easily generalized to higher or lower dimensions. The only major part that would require rework is Mesh Generation (Section 5). The process can also be implemented for time dependent systems, but we limit ourselves to time independence for our introduction.

In the process of FEM that we analyze, there are two major sections, **Theoretical**, and **Numerical**. The theoretical parts are required to be compleated prior to any numerical computation. These stages of the process, must be done manually inorder to determin the program to implement. These stages are extreamply dependent on the general form of the probelm at hand, and is where much of the mathematics

is required. Then there are the numerical computational stages. These steps are frequently independent of the problem at hand, and thus much of the work can be automated. A single implementation of the computational steps can be utilized for a general class of problems. Aswell as many of the individual steps can be used for any problem, and are compleatly independent of the rest of the process.

2. PROBLEM STATEMENT

For our expanation we will consider the 2D steady state heat equations. This is a simplification of the unsteady multidimensional heat transfer model. Since we are focusing on the time independent analysis, there is no need for an initial condition. We definine our domain to be Ω such that $\Omega \subset \mathbb{R}^2$ is a bounded 2D domain with boundary Γ .

We write the general expression for the time-dependent general heat equation as

$$\rho c_p \frac{\partial T}{\partial t} + \rho c_p \vec{u} \cdot \nabla T - \nabla \cdot (\vec{\bar{x}} \nabla T) = f \quad \text{in } \Omega.$$

Where $T(\vec{x}, t)$ is the temperature at point $\vec{x} \in \Omega$ at time $t \in (0, \tau)$. ρ is the density, and c_p denotes the specific heat capacity. For our case we assume that they are constant and 1. The forcing function f is given by $f = \rho h$ where h is the rate of heating per unit mass. \vec{u} is the velocity, and $\vec{\bar{x}}$ is the thermal conductivity is a 2×2 matrix.

Using many simplifications and our assumption that there is no change in time, we find the problem statement to be

$$(2.1) \quad \rho c_p \vec{u} \cdot \nabla T - \nabla \cdot (k \nabla T) = f \quad \text{in } \Omega.$$

Using the definition of gradient and divergence operators, we find that this simplified equation can be written as

$$\sum_{i=1}^2 u_i \frac{\partial T}{\partial x_i} - \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial}{\partial x_i} \left(k \frac{\partial T}{\partial x_j} \right) = f \quad \text{in } \Omega.$$

Note that the solution of this equations is defined with some constant, meaning that if T is a solution to this equation, then $T + C$ for $C \in \mathbb{R}$ is also a solution. In order to avoid this, we must impose some boundary conditions.

3. BOUNDARY CONDITIONS

We must prescribe boundary conditions on our problem statement (Eq 2.1), to ensure that there is a single solution. For this initial introduction we only present the process with Dirichlet boundary conditions

$$T = T_0 \quad \text{on } \Gamma,$$

where $T_0(\vec{x})$ is some function defining the temperature along each of the edges of the boundary. Note that this function can be of any form, including a piecewise function.

It is important to note that other boundary conditions, such as Neumann, and Newton boundary conditions can be used. These other boundary conditions add more changes to the formulation of the problem later on, and require extra modifications to the implementation. Because of the extra alterations that they would necessitate, we will only focus on the Dirichlet boundary conditions.

If the time-dependence of the problem had been retained, we would also require some initial condition

$$T(\vec{x}, 0) = T^0(\vec{x}).$$

that would be defined on our domain Ω . However, we are able to neglect the initial condition as we are only looking for the stationary boundary value problem solution.

4. VARIATIONAL FORMULATION

We first define the concept of a residual. The residual is the process of constructing a helper function $R(f)$, such that when f is the exact solution

$R(f) = 0$. This simplifies the process of solving this partial differential equation, into minimizing the value of the residual. This concept is used multiple times throughout the process of FEA.

We construct the residual of equation 2.1 as

$$R(\vec{x}) = \rho c_p \vec{u} \cdot \nabla T - \nabla \cdot (k \nabla T) - f.$$

Clearly $R(\vec{x}) = 0$ if and only if T is the exact solution. Multiplying the residual by some test function w and integrating, we obtain the weighted residual

$$\int_{\Omega} w(\vec{x}) R(\vec{x}) d\vec{x}.$$

For the exact solution, this will always be zero because $R(\vec{x}) = 0$. Thus, we can select any function for the test function w , and this expression can still be satisfiable. We can expand the weighted residual to obtain

$$(4.1) \quad \int_{\Omega} \rho c_p w \vec{u} \cdot \nabla T d\vec{x} - \int_{\Omega} w \nabla \cdot (k \nabla T) d\vec{x} = \int_{\Omega} w f d\vec{x}.$$

Applying Green's theorem on the second integral we are able to express it as

$$\begin{aligned} & - \int_{\Omega} w \nabla \cdot (k \nabla T) d\vec{x} \\ & = \int_{\Omega} k \nabla w \cdot \nabla T d\vec{x} - \int_{\Gamma} k w \nabla T \cdot \vec{n} ds. \end{aligned}$$

The integral over Γ goes to zero, because we require our test function to be zero on the boundary. Thus we obtain the representation of the continuous variational problem

$$(4.2) \quad \int_{\Omega} \rho c_p w \vec{u} \cdot \nabla T d\vec{x} + \int_{\Omega} k \nabla w \cdot \nabla T d\vec{x} = \int_{\Omega} w f d\vec{x}.$$

5. MESH GENERATION

For the implementation of FEM, we must discretize our global domain into a set of finite subdomains, which we call elements. We can then compute our approximation on each element independent of the rest of the domain, then combine the solutions of each element into our global solution.

These elements can be of the form of an N sided polygon. However, with polygons with more than 3 edges, issues can arise. So for the purposes of FEM, we will be focused on the construction of a mesh of triangles. This is the most common mesh polygon,

as any mesh of polygons with more edges can be represented by a triangular mesh. In addition, most research of mesh generation has been focused on triangular meshes for this same reason.

Our process will be construction what is known as a Constrained Delaunay Triangulation, and then apply a refinement algorithm to ensure that the triangles in the mesh are “nice”.

5.1. Delaunay Triangulation. Delaunay triangulation is the straight line dual of the Voronoi diagram. Delaunay triangulations are used in many different situations, including our purpose of mesh generation for finite element analysis.

Add example image of Delaunay triangulation.

A key property of Delaunay triangulation is that no point may lie inside the circumcircle of any other triangle. We use this as our definition of Delaunay Triangulation.

Definition 5.1 (Delaunay Triangulation). Let S be a set of points in the plane. A triangulation T is a *Delaunay triangulation* (DT) of S if for each triangle t of T there exists a circle C with the following properties:

- (1) the vertices of the triangle t are on the boundary of circle C
- (2) no other vertex of S is in the interior of C .

Add image showing circumcircle definition.

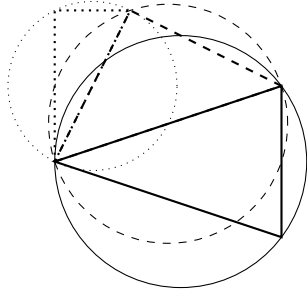


FIGURE 5.1. Demonstration of the Circumcircle definition of Delaunay.

In most situation the Delaunay triangulation of a set of points is unique. This is true for all points with one exception of a square, where there are two valid Delaunay triangulations.

A major advantage of Delaunay triangulation is that it avoids triangles with small included angles. This makes this type of triangulation extremely well suited for our purpose of FEA.

There are many different algorithms that can be used for the construction of Delaunay triangulation of a set of points in a plane. Some notable ones include;

Divide and Conquer developed by Chew[Che89], Incremental developed by Watson, Sweep Line developed by Žalik[Ž05][Dv08], and Edge Flipping developed by Sloan[Slo87][Slo93]. Each of these methods has different advantages and different efficiency in computation time.

5.2. Constrained Delaunay Triangulation. The Constrained Delaunay Triangulation(CDT), is a modification of the Delaunay Triangulation such that specified edges are forced to exist in the triangulation. This has the unfortunate effect of the resulting triangulation not being strictly Delaunay. Thus we define our constrained triangulation at the closest triangulation to the Delaunay triangulation that still includes the specified edges.

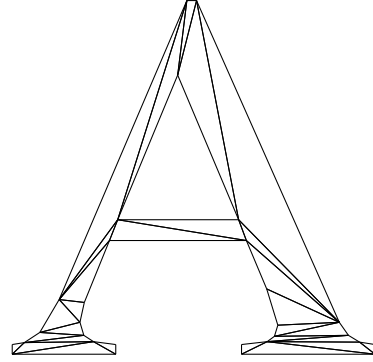


FIGURE 5.2. Constrained Delaunay triangulation.

The constrained edges are most commonly arise for the purpose of defining boundaries of the domain or enforcing a medium change. This makes it necessary for any desired domain without a strictly convex outer hull that is generated from the unconstrained triangulation.

There are several different methods of construction the CDT. A few of these generate enforce the constraints during the initial construction of the Delaunay triangulation. However, many more use a pre-constructed Delaunay triangulation and then proceed to enforce the edges, and refactor the mesh around the forced edge.

5.3. Edge Flipping. We utilize an implementation of edge flipping algorithms developed by Sloan. Sloan’s edge flipping algorithm has implementation specifics for both constrained and unconstrained Delaunay triangulation.

We will provide a short explanation of the method that is used for the construction of our triangular mesh, but more detail can be found in [Slo87][Slo93].

First Sloan's algorithm constructions the unconstrained Delaunay triangulation then enforces the required edges into the triangulation. So, we begin with the process for construction the unconstrained triangulations.

5.3.1. *Construction of Delaunay Triangulation.* Using this method also provides the advantage of generating a triangle adjacency list for each of the triangles. This allows for optimizations later in the process of finite element analysis.

There are seven main stages of the construction of the Delaunay triangulation. We also define G as the set of points to be triangulated, and N as the number of points in G .

1. Normalize coordinates of points.
2. Sort points into bins.
3. Establish the super triangle.
4. Loop over each point, repeating 5-7
5. Insert new point in triangulation.
6. Initialize stack.
7. Restore Delaunay triangulation.

In stage 1, we scale all the points that will be used to construct the triangulation between 0 and 1, this should be done such that all relative positions of points are unchanged. This is done using

$$\hat{x} = \frac{x - x_{min}}{d_{max}}, \quad \hat{y} = \frac{y - y_{min}}{d_{max}},$$

where

$$d_{max} = \max \{x_{max} - x_{min}, y_{max} - y_{min}\}.$$

This will scale all coordinates between 0 and 1 but will maintain the relative positions to all other points.

In stage 2, the normalized points are sorted into a set of bins. Each bin is associated with a rectangular portion of global space that the points are included in. We construction the bins such that each bin contains roughly \sqrt{N} points. The bins are then ordered such that adjacent bins are numbered consecutively. This is done to improve the efficiency of later searching of the triangular mesh.

In step 3, we construction a "super triangle". This triangle is a triangle of three additional points, such that all points in G are enclosed within the super triangle.

In stage 4, we repeat stages 5-7 for every P in G .

In step 5, we determine the triangle which encloses P (Since the super triangle encloses all points, then this triangle must exist for P). We now delete this triangle, and construct three new triangles, by connection P to the three vertices of the deleted triangle.

In step 6, we add up to three triangles that are adjacent to the edges opposite P to a stack (A *stack* is a last in first out data structure, like a stack of trays).

In step 7, while the stack is not empty we proceed as following

- 7.1. Remove triangle from top of the stack.
- 7.2. If P is within the circumcircle of this triangle, then the adjacent triangle containing P and the unstacked triangle form a convex quadrilateral whose diagonal is drawn in the non-optimal direction. Swap this diagonal edge.
- 7.3. Add any new triangles which are opposite P to the stack.

After the stack is empty, then Delaunay triangulation has been restored.

Once every point has been added to the triangulation, then the completed Delaunay triangulation can be returned.

5.3.2. *Construction of Constrained Delaunay Triangulation.* Now that a Delaunay triangulation has been constructed for our set of G points, we now modify the triangulation so to ensure that certain edges are present.

This process proceeds as follows

1. Loop over each constrained edge.
2. Find intersecting edges.
3. Remove intersecting edges.
4. Restore Delaunay triangulation.
5. Remove superfluous triangles.

In step 1, we loop over every edge to be constrained, and we define the edge by the endpoints of the edge V_i and V_j . Using this edge repeat steps 2-4.

In stage 2, we find all edges in the triangulation that intersect $V_i - V_j$. Store all of these edges in a list. If our constrained edge is already present in the triangulation, the go to step 1.

In step 3, repeat while there are still edges that intersect $V_i - V_j$.

- 3.1. Remove an edge from the list of intersecting edges, call this edge $V_k - V_l$.
- 3.2. If the triangles that share the edge $V_k - V_l$ do not form a strictly convex quadrilateral, then place the edge back on the list of edges and go to step 3. Else, swap this diagonal, and define the new diagonal as $V_m - V_n$. If $V_m - V_n$ still intersections $V_i - V_j$, then place it back on the list of new edges, otherwise place $V_m - V_n$ on a list of new edges.

In step 4, repeat until no changes occur.

- 4.1. Loop over each newly created edge, define each edge by $V_k - V_l$.

- 4.2. If $V_k - V_i$ is equal to $V_i - V_j$, then skip to step 4.1.
- 4.3. If the two triangles that share the edge $V_k - V_i$ do not satisfy the Delaunay criterion, then swap the diagonal, and place the new diagonal on the list of new edges.

In step 5, we need to remove all unwanted triangles. For this stage, we differ from the process of Sloan. This is because for our implementation we desire the ability to have holes in the mesh, and the simplistic method of removing exterior triangles will not achieve this.

We use what is known as a triangle infection algorithm. Where we define a few points that will infect the triangles which the points are within, then the infection spreads to all adjacent triangles that are not separated by a constrained edge. This process is repeated until no new triangles become infected. Then all infected triangles are removed.

With this process, we need to define a point that is within every one of the holes in the mesh, as the initial infection point, then the three vertices of the super triangle are also used as initial infection points.

After all infected triangles are removed, we have constructed our constrained Delaunay triangulation, which can then be used for FEA.

5.4. Mesh Refinement. The raw result of our constrained Delaunay triangulation algorithm may be sufficient in a few situations, but for most purposes, it returns sub-optimal meshes. Because of this, we must implement a mesh refinement algorithm to improve upon the quality of our mesh.

There are two main mesh refinement algorithms, Ruppert's algorithm [Rup95], and Chew's second algorithm [Che93]. These two algorithms work with a similar principle, but we will focus on Chew's second algorithm. We do this because of the advantages that Chew's second algorithm provides over Ruppert's.

Mesh refinement is required because triangular elements along constrained edges are currently forced to have the constrained edge and one of their edges. This can cause these triangles to be skinny sliver triangles, ones that are undesirable in our final mesh. To fix this we define a notion of "nice" triangles, and insert new points into the mesh until all triangles satisfy our definition of "nice".

Definition 5.2 (*Nice Triangle*). A triangle is *nice* if it is both well-shaped(Def 5.3) and well-sized(Def 5.4).

Definition 5.3 (well-shaped). A triangle is well-shaped if all its angles are greater than or equal to some angle α (commonly $\alpha = 30^\circ$).

Definition 5.4 (well-sized). A triangle is well-sized if the area of the triangle satisfies some user defined grading function g (commonly $g = \text{const}$). This function can use any criteria as long as there exists a value $\delta > 0$ such that any well-shaped(Def 5.3) triangle that fits within a circle of radius δ would satisfy the grading function.

The use of a non-constant grading function is to allow dynamically sized meshes. Such that in areas of interest there can be significantly more refinement, and areas of less interest can be generalized. The restrictions on the function simply stated is that there must be some size of a triangle that will satisfy the grading function everywhere, and that the grading function does not get infinitely strict at any point.

Using these definitions we are able to implement the mesh refinement algorithm. Chew's second algorithm proceeds as follows.

1. Grade any triangles that are currently ungraded. A triangle only passes if it is *nice*(Def 5.2).
2. If all triangles pass then Halt. Otherwise select the largest triangle that fails Δ , and determine its circumcenter c .
3. Traverse the triangulation from any vertex of Δ in the direction of c until either running into a source-edge or finding the triangle containing c .
4. If the triangle containing c was found then insert c into the triangulation, and update the triangulation to be Delaunay. This process is similar to that of sec 5.3. Then go to step 2.
5. If a source edge was encountered, then split the source edge into two equal sized edges and update the triangulation. Let l be the length of the new edges. Delete each circumcenter-vertex(vertices that are not part of the original triangulation) that is within l (line-of-sight distance, where a source-edge means that a vertex is infinitely far away) of the new vertex. Then go to step 2.

Using this algorithm it is possible to construct triangulations with a required minimum angle and force a size function. This provides the ability of guaranteed *nice* triangular meshes for any desired input domain.

Both Chew's second algorithm and Ruppert's algorithm have proven termination for minimum angles below a certain point. For Ruppert's algorithm, any minimum angle below 20.7° is guaranteed to halt. While Chew's second algorithm is guaranteed to halt for angles below 28.6° and will often succeed with angles below 34° . Because of this improvement on

the guaranteed minimum angle is why our focus is on Chew’s algorithm for mesh refinement.

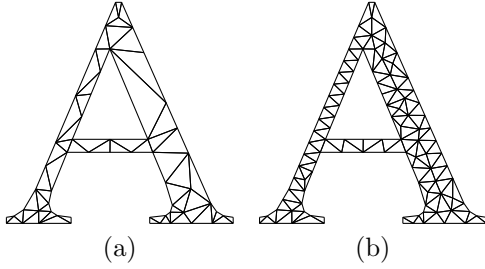


FIGURE 5.3. These four images demonstrate the refinement of a mesh. (a) 60 triangles, (b) 132 triangles, (c) 1,319 triangles, and (d) 12,999 triangles.

5.5. Implementation. There are many implementations of this element discretization process. Several of the original papers provide FORTRAN code, that can be used in implementations. We transcribed Sloan’s FORTRAN code for his implementation of Constrained Delaunay Triangulation construction into C++, which is included in [APPENDIX CODE]. However, for mesh refinement, we utilized the open source software Triangle[She96], which implements the Delaunay triangulation construction and mesh refinement.

Add code to the appendix.

6. LINEAR SYSTEMS

As the finite element method constructs a global system of equations, we need to implement a method for solving this system of linear equations. However, there are issues with the magnitude of this system. Let us take a sample situation where our mesh has $\sim 100,000$ elements. This would result in a global matrix of 10,000,000,000 values, which would require $\sim 80Gb$ of memory. Solving this matrix using straight forward gaussian elimination, which has order $\mathcal{O}(n^3)$ time complexity, would require extreme amounts of time and memory to compute. Because of these issues we need to utilize methods that allow us to avoid these issues that arise in the brute force method of solving the system of linear equations.

6.1. Data Structures. The first issue that we need to solve is the issue of memory requirements for the large matrices. We can greatly use that fact that the global matrix will be a sparse matrix. This is true by the method of construction of our matrix, most of the elements will be zero.

Using this knowledge we examine three methods of matrix storage. We use a sense of *efficiency* to imply memory efficiency, and the memory usage required to store the matrix.

We assume that the indices of the matrix are stored as `uint64_t` which requires 8 bytes of memory, and that the values are stored as `double` which requires 8 bytes of memory. We use N to denote the dimension of the square matrix, and S to represent the percent of elements that are non zero in the matrix.

1. Full matrix storage.
2. Dictionary of Keys.
3. Compressed Row Storage.

6.1.1. Full Matrix. This method simply stores a two dimensional array of `double` values. This means that the memory requirement will be $N^2 \cdot 8$ bytes. Note that this is independent of the sparsity of the matrix, so even if the matrix has only one element the memory usage is still the same as a full matrix.

An example of this is the following

$$A = \begin{bmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 0 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 0 & 0 & 5 & 0 \\ 0 & 8 & 0 & 9 & 0 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{bmatrix}$$

For this method we directly store this matrix into memory saving all of the zero elements.

This is the most simplistic method for storing a matrix. It will not suffice for most of our situations, but it is important to have a base line to compare to, to determine the efficiency of our later methods.

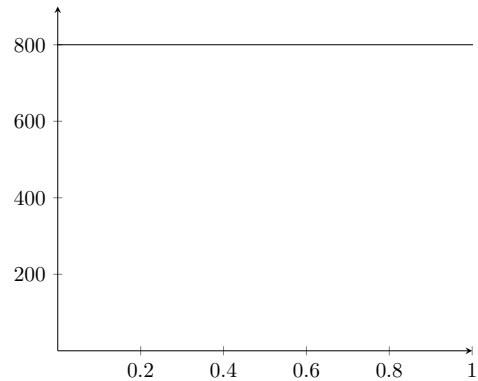


FIGURE 6.4. Memory usage for full matrix storage at different percentages of sparsity.

6.1.2. *Dictionary of Keys.* The next logical step for improving our matrix efficiency, is to only store the elements that are not zero. A simple way of doing this is to store the *row* and *column* indecies and the associated value at that point. This means that for each value we store three numbers, but we only store the values that are not zero. This means that this is only effective to an extent, as a full matrix would require three times the memory of the method in section 6.1.1.

Using the same example from above

$$A = \begin{bmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 0 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 0 & 0 & 5 & 0 \\ 0 & 8 & 0 & 9 & 0 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{bmatrix},$$

we implement the DOK method for sparse matrix storage.

$$A = \begin{bmatrix} 0 & 0 & 10 \\ 0 & 4 & -2 \\ 1 & 0 & 3 \\ 1 & 5 & 3 \\ 2 & 1 & 7 \\ 2 & 2 & 8 \\ 2 & 3 & 7 \\ 3 & 0 & 3 \\ 3 & 4 & 5 \\ 4 & 1 & 8 \\ 4 & 3 & 9 \\ 4 & 5 & 13 \\ 5 & 1 & 3 \\ 5 & 4 & 2 \\ 5 & 5 & -1 \end{bmatrix}$$

Where the first column in the row index (0 based), the second column in the column index (0 based), and the third column is the value stored at that position.

We find that this method requires $SN^2 \cdot 24$ bytes. It is clear to see that if S is large, then this method is significantly less effective but if S is small enough, then there is additional efficiency to be gained. We can compute the minimum value of S that would allow for increased efficiency like so

$$SN^2 \cdot 24 = N^2 \cdot 8$$

$$S = \frac{8}{24}$$

$$S = \frac{1}{3}.$$

Thus for any matrix where more than a third of its elements not zero, means that we would be better off using a full matrix storage system.

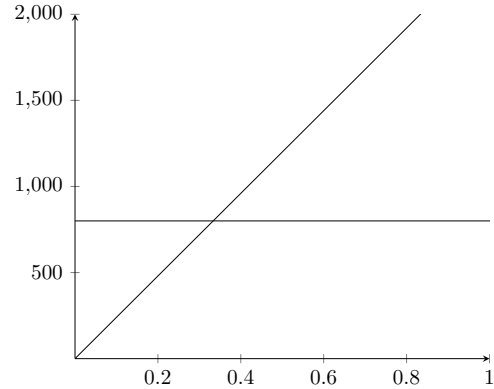


FIGURE 6.5. Memory usage for DOK at different percentages of sparsity.

6.1.3. *Compressed Row Storage.* This method take the concept of DOK, and modifies it by the realization that many elements will be on the same row, so we just need to store the number of elements in a row, the elements column index, and the value. This means that for every value, we now only need to store two numbers, and we have a list of row counts that must exists for all number of non-zero elements.

A slight optimization that we introduce, is instead of storing the number of elements in each row, we store the total number of elements so far. This is a slight improvement in the coputational implementation for element access.

This means that for every matrix, we need to store three vectors. **val**, **col_ind**, and **row_ptr**.

Using our example matrix,

$$A = \begin{bmatrix} 10 & 0 & 0 & 0 & -2 & 0 \\ 3 & 0 & 0 & 0 & 0 & 3 \\ 0 & 7 & 8 & 7 & 0 & 0 \\ 3 & 0 & 0 & 0 & 5 & 0 \\ 0 & 8 & 0 & 9 & 0 & 13 \\ 0 & 4 & 0 & 0 & 2 & -1 \end{bmatrix},$$

we find our three vectors to be

$$\text{row_ptr} = [0, 2, 4, 7, 8, 12, 15]$$

$$\text{col_ind} = [0, 4, 0, 5, 1, 2, 3, 0, 4, 1, 3, 5, 1, 4, 5]$$

$$\text{val} = [10, -2, 3, 3, 7, 8, 7, 3, 5, 8, 9, 13, 4, 2, -1]$$

Note that we can easily find the number of nonzero elements by reading the last value in **row_ptr**.

We can find the efficiency of this method of matrix storage to be $8N + SN^2 \cdot 16$. Again it is clear that

if the matrix is full, then this method is incredibly inefficient. It is also interesting to note that if the matrix is empty, then the DOK method is more efficient. There are two values of S where

1. Dictionary of keys transitions from being more efficient to less.
2. Full matrix transitions from being less efficient to more.

We solve for both of these values of S . First, for when DOK becomes less efficient than CRS.

$$\begin{aligned} 8N + SN^2 \cdot 16 &= SN^2 \cdot 24 \\ 8N &= 8SN^2 \\ S &= \frac{1}{N}. \end{aligned}$$

And when the full matrix becomes more efficient.

$$\begin{aligned} 8N + SN^2 \cdot 16 &= N^2 \cdot 8 \\ 16SN^2 &= 8N^2 - 8N \\ S &= \frac{N-1}{2N} \end{aligned}$$

It is interesting to take notice that the range in which this method is most efficient is dependent on the size of our matrix. Thus we only want to use this value when

$$\frac{1}{N} < S < \frac{N-1}{2N}.$$

Using this relation we can construct a plot for the range of S values in relation of N in which this method should be utilized.

Add filled plot of range of optimal S for given N . It looks really cool.

We can use this relation to see that for large N , this method is the most efficient between

$$0 < S < 0.5$$

Because of this asymptotic approach of this range of sparsity for maximum efficiency, we use this method of sparse matrix storage in our implementation.

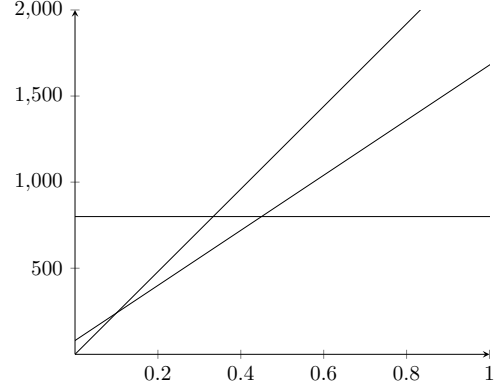


FIGURE 6.6. Memory usage for CRS and DOK at different percentages of sparsity.

6.1.4. Vectors. We will similarly have large vectors, and can implement a similar method for sparse vector storage, but this would be counter productive, as our large vectors are not likely to be sparse. Thus using a sparse representation of a vector would actually consume more memory than the full vector storage method.

6.2. Krylov Subspace. An extremely useful tool for for numerical approximation of large systems of equations is *Krylov Subspaces*. Krylov subspaces are defined by Krylov[Kry31] as

PAUL! All the sources I find about Krylov subspaces just seem to use it without any actual explanation or reasoning. So I added Krylov's original paper (Note; it is in russian, so I didn't read it), and provide the reasoning that we used. Is that sufficient?

$$\mathcal{K}_r(A, b) = \text{span} \{b, Ab, A^2b, \dots, A^{r-1}b\}$$

where b is some vector. This basis assists in the approximation for linear systems of the form $Ax = b$. The subscript r simply defines the point at which we limit our approximation, if r is infinite then the Krylov subspace is the full space, and solutions in that Krylov subspace will be exact.

Note that the basis vectors in the Krylov subspace are not necessarily orthogonal, nor normal.

6.2.1. Reasoning. We provide a short example of the reasoning for the construction of the Krylov subspace. We consider the example where we want to solve a system of linear equations of the form $Ax = b$. Where $A = I + \varepsilon$.

$$\begin{aligned}
Ax &= b \\
x &= A^{-1}b \\
x &= (I + \varepsilon)^{-1}b.
\end{aligned}$$

We construct an approximation for the inverse of a matrix like by geometric formula

$$\begin{aligned}
I + \varepsilon + \varepsilon^2 + \dots + \varepsilon^n &= S \\
\varepsilon + \varepsilon^2 + \dots + \varepsilon^{n+1} &= S\varepsilon \\
S - I + \varepsilon^{n+1} &= S\varepsilon \\
S - S\varepsilon &= I - \varepsilon^{n+1} \\
S(I - \varepsilon) &= I - \varepsilon^{n+1} \\
S &= (I - \varepsilon^{n+1})(I - \varepsilon)^{-1}.
\end{aligned}$$

When $n \rightarrow \infty$ and ε is small enough ($\|\varepsilon\| < 1$), then $\varepsilon^{n+1} \rightarrow 0$. Thus

$$\begin{aligned}
S &= (I - \varepsilon^{n+1})(I - \varepsilon)^{-1} \\
S &= (I - \varepsilon)^{-1}.
\end{aligned}$$

Finally we can conclude that

$$A^{-1} = (I + \varepsilon)^{-1} = I - \varepsilon - \varepsilon^2 - \dots$$

Using this in our expression for x we find

$$\begin{aligned}
x &= A^{-1}b = b - \varepsilon b - \varepsilon^2 b - \dots \\
&= b - (A - I)b - (A - I)^2 b - \dots \\
&= b - Ab + b - A^2 b + 2Ab - b + \dots
\end{aligned}$$

Using this we extract the basis that composes x to be

$$\{b, Ab, A^2b, A^3b, \dots\}$$

The Krylov subspace takes this infinite basis, and cuts it off at some given r value, thus providing us with

$$\{b, Ab, A^2b, \dots, A^{r-1}b\}$$

Which is our Krylov subspace. This is not a full proof of the construction of the Krylov subspace, but it should be a sufficient explanation for our purposes.

6.3. Stationary Iterative Methods. All iterative methods can be expressed in the form

$$x^{(k)} = Bx^{(k-1)} + c,$$

when B nor c are dependent on the iteration k , then this is a stationary iterative method. We comment on two stationary iterative methods *Jacobi*, and *Gauss-Seidel* methods. There are several other stationary iterative methods, but we limit ourselves to examining only these two methods.

6.3.1. Jacobi Method. The Jacobi method is a very simple method for solving linear systems. It is very easy to understand and implement, but the convergence of the approximation can be very slow.

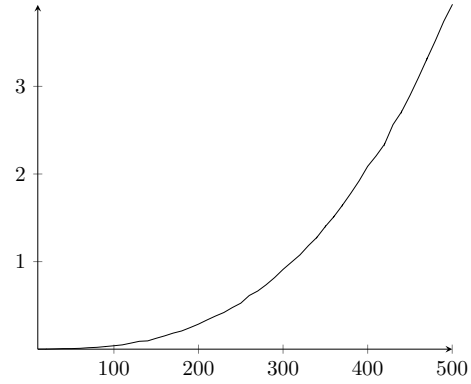


FIGURE 6.7. Time complexity of Jacobi method.

The concept for the Jacobi method, is to solve for each variable x_i one at a time, with the assumption that all other variables are constant. We solve each equation in isolation, and then the iteration provides us with our method for converging to the actual solution. We consider the i th equation in the system,

$$\sum_{j=1}^n a_{i,j}x_j = b_i.$$

We then solve for the value of x_i assuming all other values of x are fixed. We determine that

$$x_i = \frac{b_i - \sum_{j \neq i} a_{i,j}x_j}{a_{i,i}}.$$

We then implement the iteration to find the equation

$$x_i^{(k)} = \frac{b_i - \sum_{j \neq i} a_{i,j}x_j^{(k-1)}}{a_{i,i}}.$$

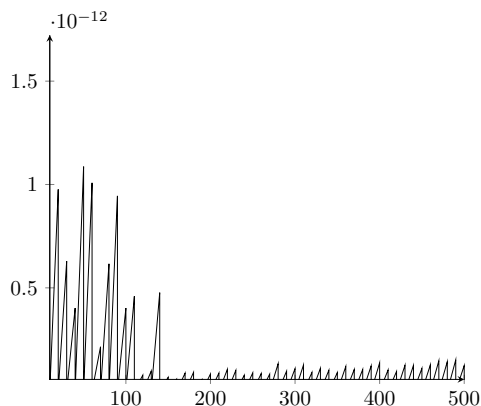


FIGURE 6.8. Error of Jacobi method.

Using this formula for the computation of the new approximation vector x

Algorithm 6.1 Jacobi Method

```

 $x = 0$ 
for  $k = 1, 2, \dots$  do
   $\bar{x} = 0$ 
  for  $i = 1, 2, \dots, n$  do
    for  $j = 1, 2, \dots, i-1, i+1, \dots, n$  do
       $\bar{x}_i += a_{i,j}x_j$ 
    end for
     $\bar{x}_i = \frac{b_i - \bar{x}_i}{a_{i,i}}$ 
  end for
   $x = \bar{x}$ 
  if  $\|Ax - b\| \leq 10^{-10}$  then
    return  $x$ 
  end if
end for
return  $x$ 

```

Pseudocode for the Jacobi method is shown in Algorithm 6.1. Note the necessity for a temporary \bar{x} term, as the values of x are needed for the computation of the next iteration. We stop the iteration either when a maximum iteration count is reached, or when our approximation is within 10^{-10} of the actual solution. It is possible to alter this tolerance to achieve either more precise approximation, or to accelerate convergence at the cost of accuracy.

6.3.2. Gauss-Seidel Method. The Gauss-Seidel method is an extension of the Jacobi method, with one optimization. Instead of using the old values of x , the algorithm uses the already computed values of the current iteration.

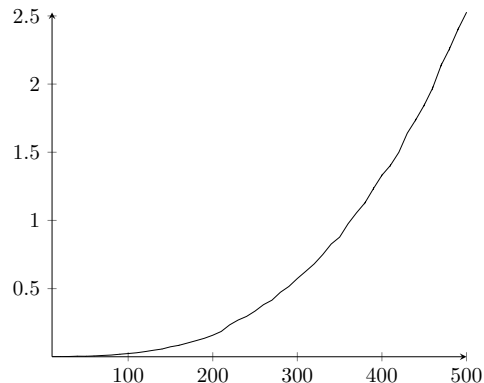


FIGURE 6.9. Time complexity of Gauss-Seidel method.

This method does not improve the rate of convergence significantly, but it is a simple to implement improvement over the Jacobi method. The mathematics and the pseudocode are almost identical to that of the Jacobi method, with the only difference being that instead of using \bar{x} all values are taken from x .

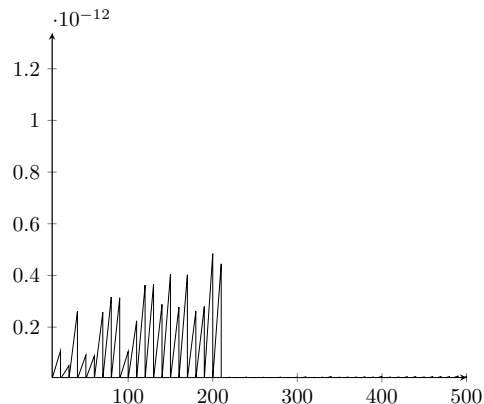


FIGURE 6.10. Error of Gauss-Seidel method.

Because of this similarity we do not explain in detail the mathematics, but the pseudocode is provided in Algorithm 6.2.

Algorithm 6.2 Gauss-Seidel Method

```

 $x = 0$ 
for  $k = 1, 2, \dots$  do
  for  $i = 1, 2, \dots, n$  do
     $x_i = 0$ 
    for  $j = 1, 2, \dots, i-1, i+1, \dots, n$  do
       $x_i += a_{i,j}x_j$ 
    end for
     $x_i = \frac{b_i - x_i}{a_{i,i}}$ 
  end for
  if  $\|Ax - b\| \leq 10^{-10}$  then
    return  $x$ 
  end if
end for
return  $x$ 

```

6.4. Nonstationary Iterative Methods. Nonstationary methods differ from their stationary counterparts, because the computation involve information that is altered from every iteration. E.g. the values of B and c can depend on the iteration count k . We discuss the nonstationary iterative method called the *Conjugate Gradient* method. All of the nonstationary iterative methods utilize Krylov subspaces, with the exception of a few which we do not discuss.

6.4.1. Conjugate Gradient Method. This method is extremely effective for symmetric positive definite systems, and it is one of the most researched nonstationary methods. This method is based on the construction of a sequence of approximations. Each new approximation is generated based on the direction required to minimize the residual.

Add plot of time complexity and error.

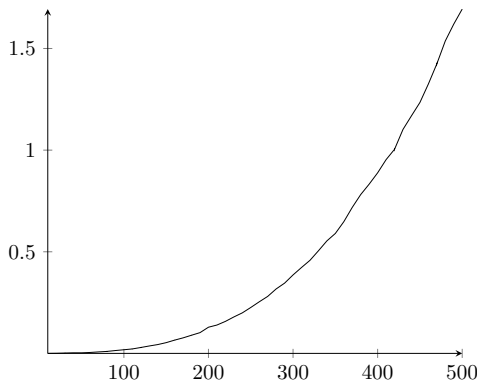


FIGURE 6.11. Time complexity of Conjugate Gradient method.

This is done by determining which “direction” must be used in order to minimize the residual, then the approximation is moved by some α amount in that direction. After repeating this process for some

number of iterations, the generated approximation should be converging to the actual solution.

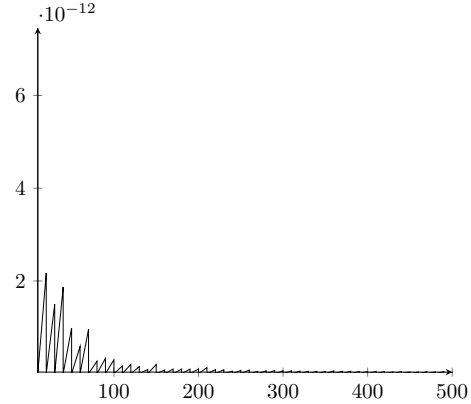


FIGURE 6.12. Error of Conjugate Gradient method.

We construct the iterative approximation in each iteration by a multiple of α_i of the search direction $p^{(i)}$,

$$x^{(i)} = x^{(i-1)} + \alpha_i p^{(i)}.$$

We then update the residual $r^{(i)} = b - Ax^{(i)}$ as

$$\begin{aligned} q^{(i)} &= Ap^{(i)} \\ r^{(i)} &= r^{(i-1)} - \alpha_i q^{(i)} \end{aligned}$$

We construct α_i in order to minimize $r^{(i)T} A^{-1} r^{(i)}$,

$$\alpha_i = \frac{r^{(i-1)T} r^{(i-1)}}{p^{(i)T} A p^{(i)}}.$$

Then to update the search directions using the residual

$$\begin{aligned} \beta_i &= \frac{r^{(i)T} r^{(i)}}{r^{(i-1)T} r^{(i-1)}} \\ p^{(i)} &= r^{(i)} + \beta_{i-1} p^{(i-1)} \end{aligned}$$

where the choice of β_i ensures that $r^{(i)}$ and $r^{(i-1)}$ are orthogonal.

This process ties together the construction of the Krylov subspace, and the construction of the approximation based off of the most recently computed Krylov basis vector. This means that if we allow our iterations to proceed until N iterations, then we should have constructed the complete Krylov subspace, and thus our approximation should be exact.

Using this approximation we can construct our pseudocode of the algorithm. Which is presented in Algorithm 6.3.

Algorithm 6.3 ConjugateGradient

```

 $x = 0$ 
 $r = b - Ax$ 
for  $i = 1, 2, \dots, n$  do
   $\rho^{(i)} = \|r\|^2$ 
  if  $i = 0$  then
     $p = r$ 
  else
     $\beta = \frac{\rho^{(i)}}{\rho^{(i-1)}}$ 
     $p = r + \beta p$ 
  end if
   $q = Ap$ 
   $\alpha = \frac{\rho^{(i)}}{\langle p, q \rangle}$ 
   $x += \alpha p$ 
   $r -= \alpha q$ 
  if  $\|r\| \leq 10^{-10}$  then
    return  $x$ 
  end if
end for
return  $x$ 

```

6.5. Cholesky Method. The Cholesky method is a method for solving systems of linear equations based upon the Cholesky decomposition of a matrix. This is the only direct method that we consider. A direct method does not depend on convergence of a sequence, but instead directly solves the system of equations. This method is useful for the purpose of time prediction. The convergence of iterative methods greatly depends on the system being considered, using this method the time required is directly related to the number of elements. This means that each system of the same size will require the same time to solve.

Add Plot of time complexity and error for N .

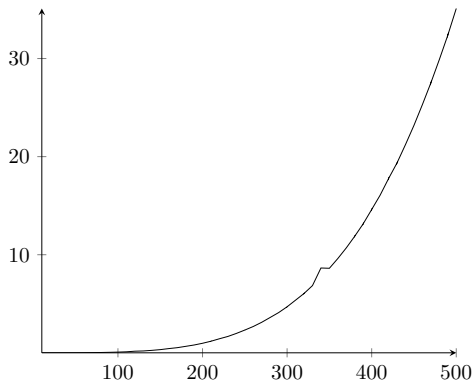


FIGURE 6.13. Time complexity of Cholesky method.

The Cholesky decomposition of a matrix provides a lower triangular matrix L of the form

$$A = LL^T$$

We use this decomposition of A to solve two systems of linear equations, using forward and back substitution, which is an efficient method for finding these solutions.

$$Ly = b \quad \text{by forward substitution}$$

$$L^T x = y \quad \text{by back substitution.}$$

This provides the solution x to the system of linear equations.

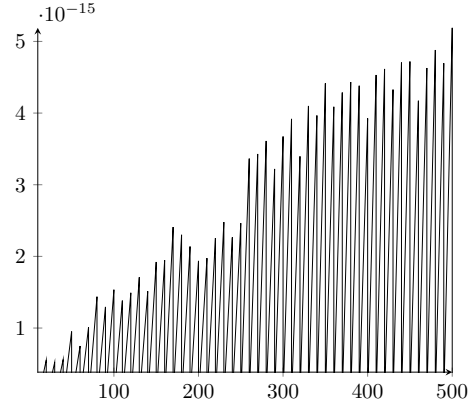


FIGURE 6.14. Error of Cholesky method.

We present a plot of the time requirements for the computation in Figure 6.13, and error in Figure 6.14. Note that the error spikes are due to floating point arithmetic error. But all of these errors are unavoidable, and should not effect any solutions too much.

6.5.1. Cholesky Decomposition. We utilize the Cholesky-Banachiewicz algorithm for the construction of the lower triangular matrix L . This algorithm states that

$$L_{j,j} = \sqrt{A_{j,j} - \sum_{k=1}^{j-1} L_{j,k}^2}$$

$$L_{i,j} = \frac{1}{L_{j,j}} \left(A_{i,j} - \sum_{k=1}^{j-1} L_{i,k} L_{j,k} \right) \quad \text{for } i > j.$$

This construction of the Cholesky decomposition is very easy, and fairly efficient. Thus allowing us to preform this decomposition easily.

6.5.2. *Forward/Back Substitution.* The method of forward and back substitution is an efficient method for solving systems of linear equations of the form $Lx = b$ where L is a lower triangular matrix. Note that the process is almost identical for an upper triangular matrix, just with the order of iteration flipped. We will only provide an explanation for forward substitution.

The first stage is to write the system of equations like so

$$\begin{array}{ccccccc} L_{1,1}x_1 & & & & & & = b_1 \\ L_{2,1}x_1 & +L_{2,2}x_2 & & & & & = b_2 \\ \vdots & \vdots & \ddots & & & \vdots & \\ L_{n,1}x_1 & +L_{n,2}x_2 & \cdots & L_{n,n}x_n & = & b_n \end{array}$$

We can clearly solve for x_1 directly, and find

$$x_1 = \frac{b_1}{L_{1,1}}.$$

It is possible to substitute this into the second equation. This will cause the second equation to only have one unknown x_2 , and so we can directly solve for x_2 . This process is repeated for all x . Providing the general expression

$$x_m = \frac{b_m - \sum_{i=1}^{m-1} L_{m,i}x_i}{L_{m,m}}.$$

This expression can then be used to compute the solution to the system of linear equations.

Using the Cholesky decomposition and forward and back substitution, we are able to compute the exact solution to the system of equations, without resorting to the inefficient Gaussian elimination.

We provide pseudocode of this method in Algorithm 6.4.

Algorithm 6.4 Cholesky method

```

    ▷ Calculate  $A = LL^T$  decomposition.
for  $j = 0, 1, 2, \dots, N$  do
     $L_{j,j} = 0$ 
    for  $k = 0, 1, 2, \dots, j$  do
         $L_{j,j} += L_{j,k}^2$ 
    end for
     $L_{j,j} = \sqrt{A_{j,j} - L_{j,j}}$ 
    for  $i = j, j+1, j+2, \dots, N$  do
         $L_{i,j} = 0$ 
        for  $k = 0, 1, 2, \dots, j$  do
             $L_{i,j} += L_{i,k}L_{j,k}$ 
        end for
         $L_{i,j} = \frac{A_{i,j} - L_{i,j}}{L_{j,j}}$ 
    end for
end for

    ▷ Calculate  $Ly = b$  by forward substitution.
 $y = 0$ 
for  $i = 1, 2, \dots, N$  do
     $y_i = 0$ 
    for  $j = 1, 2, \dots, i$  do
         $y_i += L_{i,j}y_j$ 
    end for
     $y_i = \frac{b_i - y_i}{L_{i,i}}$ 
end for

    ▷ Calculate  $L^T x = y$  by back substitution.
 $x = 0$ 
for  $i = N, N-1, \dots, 1$  do
     $x_i = 0$ 
    for  $j = i+1, i+2, \dots, N$  do
         $x_i += L_{i,j}^T x_j$ 
    end for
     $x_i = \frac{y_i - x_i}{L_{i,i}^T}$ 
end for
return  $x$ 

```

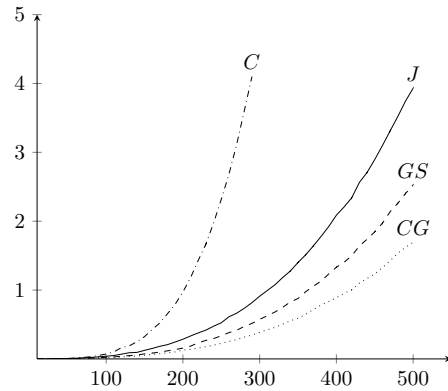


FIGURE 6.15. Plots of the different time requirements for solving the linear systems.

6.6. Issues. We note that these algorithms have requirements on the form of the matrix A . For many of these the algorithm requires a symmetric positive definite matrix. We comment on some of these restrictions here.

We note that future work must be done to determine if the global matrix that is constructed by the Finite Element Method will consistently satisfy these restraints. If it does then it is possible to provide extra optimizations based on the form of the global matrix. If a general form of the matrix cannot be guaranteed, then we will be forced to utilize significantly slower methods of computing the solution to the linear system.

6.6.1. Positive Definite. All of these algorithms require a positive definite matrix. This is required to guarantee that there is a unique solution to the system. This can be considered to ensure that there is a minimum in some N dimensional space, instead of a saddle, or maximum. Only with a minimum can we minimize our residual and use gradient decent to find our equilibrium solution.

6.6.2. Symmetric. Most of these algorithms, with the exception of Cholesky method, require the matrix to be symmetric. This requirement is mostly for computational efficiency, and simplicity of the algorithm. There are some methods that can be adapted to function for nonsymmetric matrices, but those become significantly more complex. For most purposes, we should be able to make use of our implementations of these algorithms for finding the solution to the linear system.

REFERENCES

- [BBC⁺93] Richard Barrett, Michael Berry, Tony F. Chan, James Demmel, June M. Donato, Jack Dongarra, Victor Eijkhout, Roldan Pozo, Charels Romine, and Hen Van der Vorst. *Templates for the Solution of Linear Systems: Buildign Blocks for Iterative Methods*. Society for Industrial and Applied Mathematics, 2 edition, 1993.
- [Che89] L. Paul Chew. Constrained delaunay triangulations. *Algorithmica*, 4:92–108, 1989.
- [Che93] L. Paul Chew. Guaranteed-quality mesh generation for curved surfaces. *9th Annual Symposium on Computation Geometry*, pages 275–280, 1993.
- [Dv08] V. Domiter and B. Žalik. Sweep-line algorithm for constrained delaunay triangulation. *International Journal of Geographical Information Science*, 22(4):449–462, 2008.
- [KH15] Dmitri Kuzmin and Jari Hämäläinen. *Finite Element Methods for Computational Fluid Dynamics*. Society for Industrial and Applied Mathematics, Philadelphia, 2015.
- [Kry31] A. N. Krylov. On the numerical solution of the equation by which in technical questions frequencies of small oscillations of material systems are determined. *News of Academy of Sciences of the USSR*, 7(4):491–539, 1931.
- [Rup95] Jim Ruppert. A delaunay refinement algorithm for quality 2-dimensional mesh generation. *Journal of Algorithms*, 18(3):548–585, 1995.
- [She96] Jonathan Richard Shewchuk. Triangle: Engineering a 2D Quality Mesh Generator and Delaunay Triangulator. In Ming C. Lin and Dinesh Manocha, editors, *Applied Computational Geometry: Towards Geometric Engineering*, volume 1148 of *Lecture Notes in Computer Science*, pages 203–222. Springer-Verlag, May 1996. From the First ACM Workshop on Applied Computational Geometry.
- [She02] Jonathan Richard Shewchuk. Delaunay refinement algorithms for triangular mesh generation. *Computational Geometry: Theory and Applications*, 22(1-3):21,74, 2002.
- [Slo87] S. W. Sloan. A fast algorithm for construction delaunay triangulations in the plane. *Adv. Eng. Software*, 9(1):34–42, 1987.
- [Slo93] S. W. Sloan. A fast algorithm for generating constrained delaunay triangulations. *Computers & Structures*, 47(3):441–450, 1993.
- [SS86] Youcef Saad and Martin H. Schultz. Gmres: A generalized minimal residual algorithm for solving nonsymmetric linear systems. *Sci. Stat. Comput.*, 7(3):856–869, 1986.
- [Ž05] Borut Žalik. An efficient sweep-line delaunay triangulation algorithm. *Computer-Aided Design*, 37:1027–1038, 2005.