

CE394M: FEM solvers and errors

Krishna Kumar

University of Texas at Austin

krishnak@utexas.edu

February 11, 2019

Overview

1 Solving non-linear problems

- Iterative solvers
- Tangent stiffness
- Newton Raphson

2 Arc length method

3 Error estimates

- A priori estimates
- Posterior error estimation

4 Time-dependent problems

- Explicit schemes
- Implicit method

Linear and non-linear problems

- **Linear problems**

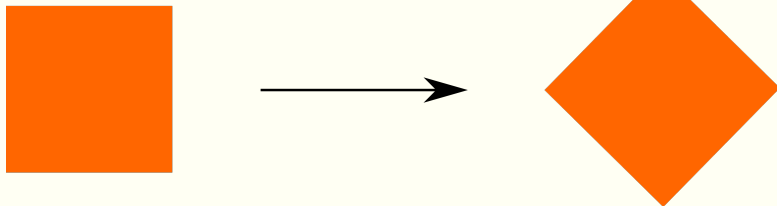
- The response can only be approximated as linear if its deformations/motions are small.
- In linear analyses, the response to individual load cases can be scaled and added to the results from other linear analyses, which is the principle of superposition.

- **Non-linear problems**

- Superposition is invalid.
- The solution is an incremental/iterative process.
- An iteration is the solution of a system of equations linearised about the current state of the nonlinear physical problem.

Non-linearity in geotechnical engineering

- **Material non-linearity**
 - plasticity
- **Contact**
 - discontinuous source of non-linearity
- **Large deformations and motions of a geotechnical structures**
 - rotation, rigid body motion
 - often ignored, still in the research area.



No strains if linear strain-displacement relations are used.

Linear solvers

To solve a system of linear equations of the form $\mathbf{K}\mathbf{a} = \mathbf{b}$, there are two families of methods of solvers that can be used: direct and iterative.

- Direct solvers solve a system of linear equations in a predefined number of steps.
- Methods are based on Gauss elimination, with the most common method being LU decomposition.
- The time required to solve a linear system increases with the number of computer operations performed.
- For a direct linear solver applied to a dense system of size $n \times n$: CPU time $= Cn^3$.
- If the number of degrees of freedom in a finite element simulation is doubled, the time required to solve the system will increase by a factor of 8!

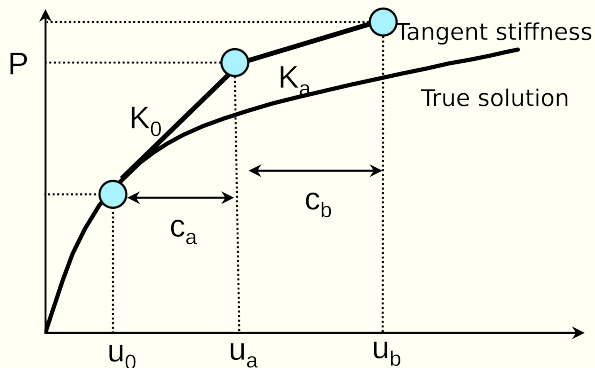
The system of equations that we derived from the finite element approximation of the BVP:

$$F_{\text{int}}(u) - F_{\text{ext}} = 0$$

At this point we remind ourselves that in the case of finite deformations, both F_{int} and F_{ext} are in general very nonlinear terms.

The integration has to be carried over the current volume and surface (of the finite element under consideration) that may in general depend on u in a highly non-linear fashion.

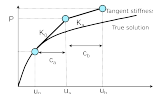
Tangent stiffness method



- Many small increments are needed to obtain accurate solution
- Need to perform a parametric study to find the optimum incremental
- Defining increments is very important
- Program - SAGE CRISP

CE394M: solvers - errors

- └ Solving non-linear problems
 - └ Tangent stiffness
 - └ Tangent stiffness method

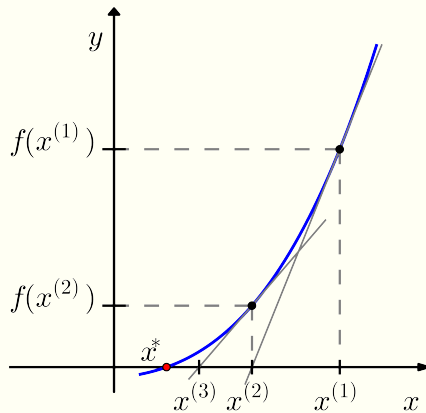


- Many small increments are needed to obtain accurate solution
- Need to perform a parametric study to find the optimum incremental
- Defining increments is very important
- Program - SAGE CRISP

In the incremental solution we divide the load into increments $\Delta P = \lambda P$, where λ is also known as a load factor and apply a repeated solution of $\Delta u = K^{-1} \Delta P$.

Basically we divide the load into substeps, and treat each as linear - but that is usually not accurate enough and inefficient.

Newton Raphson method



$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Newton Raphson method

The incremental loading is expressed as follows. The external load vector F_{ext} is gradually increased from 0 in order to reach a desired value \mathbf{F}^* . Assuming that \mathbf{F}^* itself remains constant during the analysis in terms of its 'direction' and only its magnitude is changing, we can write $F_{\text{ext}} = q$ known just to simplify our expression for the system of equations. Then we can control how the external load vector increases or decreases by introducing a scalar quantity λ and express the system as follows:

$$R(u) = F_{\text{int}} - F_{\text{ext}} \rightarrow R(u) = F_{\text{int}} - \lambda F_{\text{ext}} = 0$$

Thus by increasing or decreasing λ we can control our load. We are interested in u and λ . At every increment, we change slightly the value of λ and try to determine u satisfying $R(u) = 0$.

$$\Delta u = [K_T]_{u0}^{-1} \cdot (\Delta \lambda q)$$

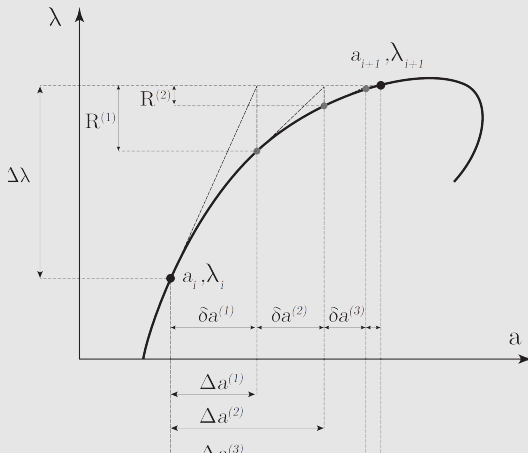
The incremental loading is expressed as follows. The external load vector F_{ext} is gradually increased from 0 in order to reach a desired value F^* . Assuming that F^* itself remains constant during the analysis in terms of its 'direction' and only its magnitude is changing, we can write $F_{ext} = q$ known just to simplify our expression for the system of equations. Then we can control how the external load vector increases or decreases by introducing a scalar quantity λ and express the system as follows:

$$R(u) = F_{int} - F_{ext} \rightarrow R(u) = F_{int} - \lambda F_{ext} = 0$$

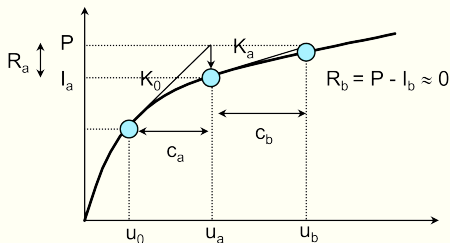
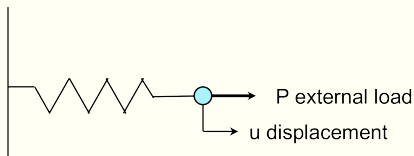
Thus by increasing or decreasing λ we can control our load. We are interested in u and λ . At every increment, we change slightly the value of λ and try to determine u satisfying $R(u) = 0$.

$$\Delta u = [K_T]_{10}^{-1} \cdot (\Delta \lambda q)$$

where $[K^T] = [\partial F(u)/\partial u]$ is the "Jacobian" matrix of the system of equations and is commonly referred to as the Stiffness Matrix.

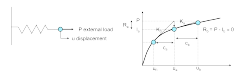


Newton Raphson method



- Using the initial stiffness K_0 , apply an increment of load δP , calculate an approximate solution c_a caused by this increment.
- The stiffness K_a is updated using the new position, and the internal force in the spring I_a is calculated.
- If the difference R_a between the total load applied to the spring, P , and I_a is smaller than the tolerance, $u_a = u_0 + c_a$ is the converged solution.

- └ Solving non-linear problems
 - └ Newton Raphson
 - └ Newton Raphson method



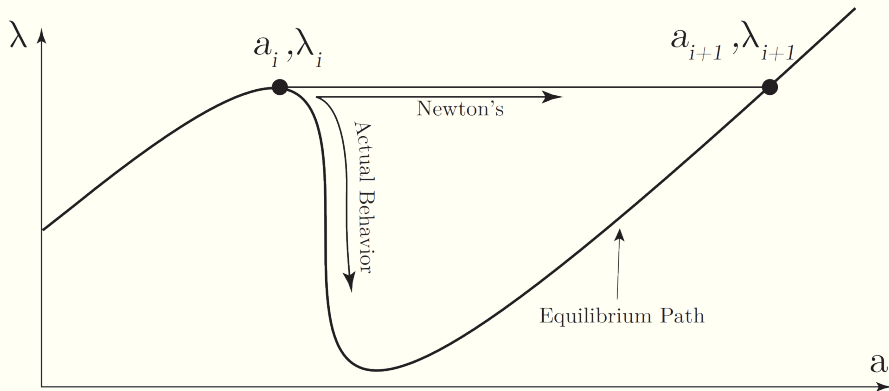
- Using the initial stiffness K_0 , apply an increment of load δP , calculate an approximate solution c_0 caused by this increment.
- The stiffness K_0 is updated using the new position, and the internal force in the spring I_0 is calculated.
- If the difference R_0 between the total load applied to the spring, P , and I_0 is smaller than the tolerance, $u_0 = c_0 + c_0$ is the converged solution.

- If R_a is not small, a new displacement correction c_b is calculated by solving $c_b = R_a/K_a$
- The new displacement u_b is updated, and the internal force I_b in the updated configuration is calculated.
- The new force residual R_b is obtained. If $R_b < \text{tolerance}$, the solution is converged. If not, continue the iteration.

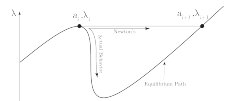
Newton Raphson method: Tolerance and Convergence

- Program - ABAQUS
- Newton-Raphson is the most standard method to solve nonlinear problems in FE.
- A large error in the initial estimate can contribute to non-convergence of the algorithm.
- **Tolerance**
 - must be small enough to ensure that the approximate solution is close to the exact mathematical solution.
 - must be large enough so that reasonable number of iterations are performed.
- **Quadratic convergence**
 - If the tangent stiffness is calculated correctly, R should reduce quadratically from one iteration to the next.

Newton Raphson: Drawbacks

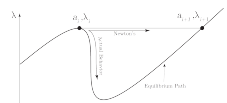


- └ Solving non-linear problems
 - └ Newton Raphson
 - └ Newton Raphson: Drawbacks



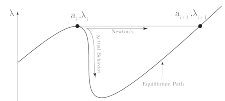
Newton's method fails to accurately follow the 'equilibrium' path once the tangent stiffness reaches zero. That happens due to the formulation of Newton's method, and in particular that it restricts the parameter λ to change monotonically every increment 3. The definition of a limit point then (saddle points excluded), suggests that in order to remain on the equilibrium path you need to change your loading pattern depending on whether the limit point is a local maximum or maximum in the $u - \lambda$ space.

- └ Solving non-linear problems
 - └ Newton Raphson
 - └ Newton Raphson: Drawbacks



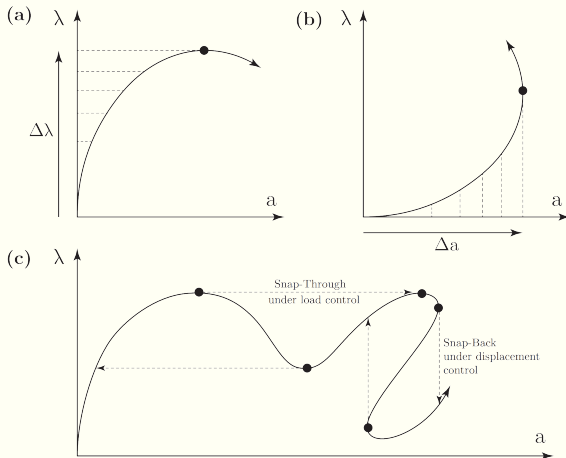
In terms of mechanical systems then, this method is able to solve any non-linear system of equations very efficiently but only up to the critical point (if any). In the case shown in figure, Newton's method fails in load-control. Now in many cases, one way to circumvent problems like these is to use displacement control, where you can continuously increase the displacements u and still remain on the equilibrium curve. In general however, apart from Snap-Through behaviors under load control, a problem may exhibit Snap-Back behaviors under displacement control or even both. The main problem is, that in most actual applications, the structural response, and therefore the equilibrium path, for the structure under consideration is unknown, and therefore one does not know what type of behavior to expect.

- └ Solving non-linear problems
 - └ Newton Raphson
 - └ Newton Raphson: Drawbacks



As a general rule, if the problem under consideration requires information after its critical/failure points then Newton's method is not a good choice. Buckling analysis and non-linear materials that exhibit work softening are just two example problems that cannot be solved using Newton's method. Furthermore, very often, strong nonlinearities that arise in finite deformation problems may eventually lead to such behaviors and thus it is necessary to introduce a numerical technique to solve such problem with strongly nonlinear behaviors.

Newton Raphson: Drawbacks



(a) A system that is unstable under load control (Snap-Through instability), (b) A system that is unstable under displacement control (Snap-Back Instability), (c) A system that is unstable under both displacement and load control

Arc length method

System of the non-linear (in general) equations to solve:

$$F_{\text{int}} - F_{\text{ext}} = 0 \rightarrow F_{\text{int}} - \lambda q = 0$$

Suppose (u_0, λ_0) satisfy the system of equations and thus belongs to the 'equilibrium path'. Arc Length postulates a simultaneous variation in both the displacements Δu and the load vector coefficient $\Delta \lambda$.

$$R(u', \lambda') = F_{\text{int}}(u_0 + \Delta u) - (\lambda_0 + \Delta \lambda)q = 0$$

If the above is satisfied for $(u_0 + \Delta u, \lambda_0 + \Delta \lambda)$ then this point also belongs to the 'equilibrium path' and we can successfully update the solution. However, immediate satisfaction is not possible, so we need to provide the necessary corrections $(\delta u, \delta \lambda)$, aiming that the new point $(u_0 + \Delta u + \delta u, \lambda_0 + \Delta \lambda + \delta \lambda)$ will satisfy.

$$R(u'', \lambda'') = F_{\text{int}}(u_0 + \Delta u + \delta u) - (\lambda_0 + \Delta \lambda + \delta \lambda)q = 0$$

CE394M: solvers - errors

└ Arc length method

└ Arc length method

Arc length method

System of the non-linear (in general) equations to solve:

$$F_{int} - F_{ext} = 0 \rightarrow F_{int} - \lambda q = 0$$

Suppose (u_0, λ_0) satisfy the system of equations and thus belongs to the 'equilibrium path'. Arc Length postulates a simultaneous variation in both the displacements Δu and the load vector coefficient $\Delta \lambda$.

$$R(u', \lambda') = F_{int}(u_0 + \Delta u) - (\lambda_0 + \Delta \lambda)q = 0$$

If the above is satisfied for $(u_0 + \Delta u, \lambda_0 + \Delta \lambda)$ then this point also belongs to the 'equilibrium path' and we can successfully update the solution. However, immediate satisfaction is not possible, so we need to provide the necessary corrections $(\delta u, \delta \lambda)$, aiming that the new point $(u_0 + \Delta u + \delta u, \lambda_0 + \Delta \lambda + \delta \lambda)$ will satisfy:

$$R(u'', \lambda'') = F_{int}(u_0 + \Delta u + \delta u) - (\lambda_0 + \Delta \lambda + \delta \lambda)q = 0$$

Arc length is a very efficient method in solving non-linear systems of equations when the problem under consideration exhibits one or more critical points. In terms of a simple mechanical loading-unloading problem, a critical point could be interpreted as the point at which the loaded body cannot support an increase of the external forces and an instability occurs.

Unlike the Newton-Method, the Arc Length method postulates a simultaneous variation in both the displacements Δu and the load vector coefficient $\Delta \lambda$. The main difference is that both Δu and $\Delta \lambda$ are unknowns in contrast to Newton's method where $\Delta \lambda$ was given and we had to iteratively solve for Δu .

CE394M: solvers - errors

└ Arc length method

└ Arc length method

Arc length method

System of the non-linear (in general) equations to solve:

$$F_{int} - F_{ext} = 0 \rightarrow F_{int} - \lambda q = 0$$

Suppose (u_0, λ_0) satisfy the system of equations and thus belongs to the 'equilibrium path'. Arc Length postulates a simultaneous variation in both the displacements Δu and the load vector coefficient $\Delta \lambda$.

$$R(u', \lambda') = F_{int}(u_0 + \Delta u) - (\lambda_0 + \Delta \lambda)q = 0$$

If the above is satisfied for $(u_0 + \Delta u, \lambda_0 + \Delta \lambda)$ then this point also belongs to the 'equilibrium path' and we can successfully update the solution. However, immediate satisfaction is not possible, so we need to provide the necessary corrections $(\delta u, \delta \lambda)$, aiming that the new point $(u_0 + \Delta u + \delta u, \lambda_0 + \Delta \lambda + \delta \lambda)$ will satisfy:

$$R(u', \lambda') = F_{int}(u_0 + \Delta u + \delta u) - (\lambda_0 + \Delta \lambda + \delta \lambda)q = 0$$

The system of equations takes the form:

$$[K^T]_{u_0 + \Delta u} \cdot \delta u - \delta \lambda q = -[F_{int}(u_0 + \Delta u) - (\lambda_0 + \Delta \lambda)q] = -R(u', \lambda')$$

δu and $\delta \lambda$ are the unknowns for whom we need to solve. If the u vector however, has dimensions $N \times 1$ then we have a total of N equations that we need to solve for $N + 1$ unknowns (N unknowns δu and 1 unknown $\delta \lambda$). Equations then are not sufficient to determine $\delta u, \delta \lambda$. The supplementary equation that completes the system is called the Arc Length Equation.

Arc length method

The Arc length equation:

$$(\Delta \mathbf{u} + \delta \mathbf{u})^T \cdot (\Delta \mathbf{u} + \delta \mathbf{u}) + \psi^2 (\Delta \lambda + \delta \lambda)^2 (\mathbf{q}^T \cdot \mathbf{q}) = \Delta l^2$$

where ψ and Δl are user defined parameters. In a sense Δl defines how far to search for the next equilibrium point and it is analogous (but not directly equivalent) to the load increment $\Delta \lambda$ we used in Newton's method.

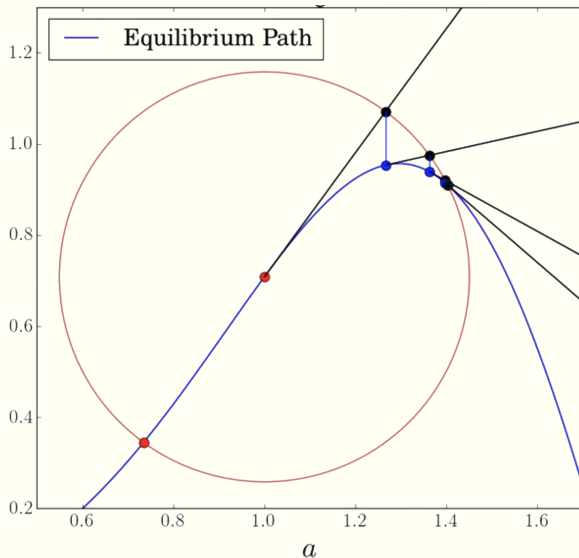
The Arc length equation:

$$(\Delta \mathbf{u} + \delta \mathbf{u})^T \cdot (\Delta \mathbf{u} + \delta \mathbf{u}) + \psi^2 (\Delta \lambda + \delta \lambda)^2 (\mathbf{q}^T \cdot \mathbf{q}) = \Delta l^2$$

where ψ and Δl are user defined parameters. In a sense Δl defines how far to search for the next equilibrium point and it is analogous (but not directly equivalent) to the load increment $\Delta \lambda$ we used in Newton's method.

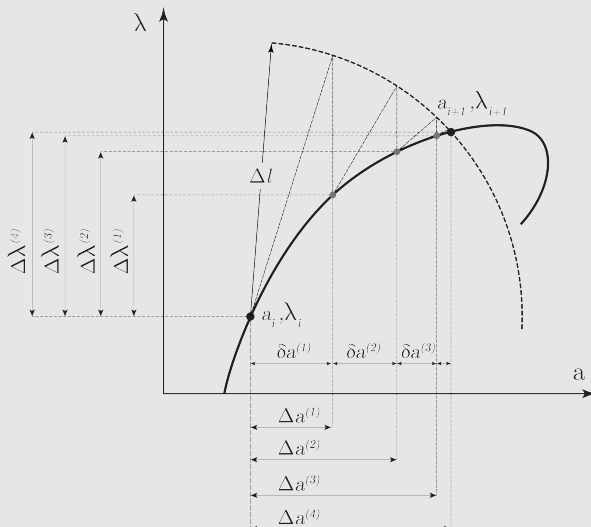
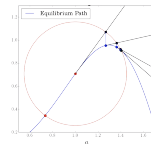
The system of equations is solved for $\delta u, \delta \lambda$ and updates the previous corrections $\Delta u, \Delta \lambda$ to be $\Delta u' = \Delta u + \delta u$ and $\Delta \lambda' = \Delta \lambda + \delta \lambda$ respectively. The method continues to provide such incremental corrections $\delta u, \delta \lambda$ until convergence is achieved. When $\psi = 1$, the method is also called the **Spherical Arc-Length Method** because the points $\Delta u', \Delta \lambda'$ belong to a circle with radius Δl . In its most general form for arbitrary ψ can be geometrically interpreted as a hyper-ellipse in the multidimensional displacement-load space $(u - \lambda)$.

Arc length method

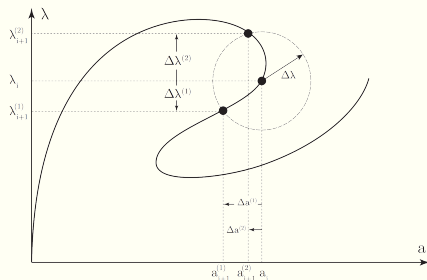
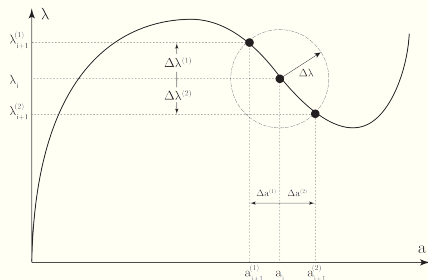


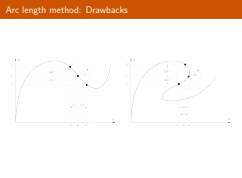
└ Arc length method

└ Arc length method

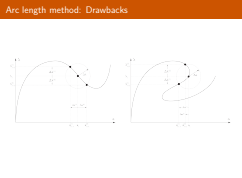


Arc length method: Drawbacks





Crisfield's [1983] implementation of Arc-length however, leads to one of the method's most important drawbacks. The quadratic equation would in general lead to two distinct solutions for $\delta\lambda$ which will in turn lead to two distinct solutions for δu . Thus, every iteration, the solver determined two sets of solutions, namely $(\delta u_1, \delta\lambda_1)$ and $(\delta u_2, \delta\lambda_2)$. This is no surprise since a circle (or a hyperellipse) would always intersect a curve in two points if its center is located on the curve.



The issue that arises then, is to develop a robust algorithm that would be able to accurately determine the correct set of $(\delta u, \delta \lambda)$ to update the solution. In general, we would like to choose the set, so that the solution 'evolves forwards'. This term 'forward evolution' is commonly used in the context of this method since choosing the wrong set would make the solution move back towards a previously converged point, and not in the desired (forward) direction. It is interesting to note, that an effective solution to this problem that works for all applications is yet to be found and as a result, many times programs like ABAQUS fail to converge to the correct solution or fail to 'evolve forwards'.

A priori estimates

- A priori error estimation is the analysis of the errors in a finite element simulation before actually performing an analysis.
- It is an abstract concept since it cannot quantify the error in a computed simulation, but analysis can tell something about whether the finite element solution will converge to the exact solution as the mesh is refined, and how quickly it will converge if the element type is changed or the mesh is refined.

Functional norms

Before considering the error, we need to decide how it should be measured. The question is how to measure a function? The answer to this is to use function norms. A common norm is known as the L^2 norm. For a function on a domain Ω , the L^2 norm is defined as:

$$\|u\|_0 = \left(\int_{\Omega} u^2 d\Omega \right)^{1/2}$$

This norm provides a measure of 'how big' a function is. A measure of the difference between two functions u and v is provided by:

$$\|u - v\|_0 = \left(\int_{\Omega} (u - v)^2 d\Omega \right)^{1/2}$$

It is only possible that $\|u - v\|_0 = 0$ if $u = v$.

Functional norms

The L^2 norm measures the 'magnitude' of a function, but does not reflect any details of the derivatives. Other norms are the H^1 norm:

$$\|u\|_0 = \left(\int_{\Omega} u^2 + \nabla u \cdot \nabla u d\Omega \right)^{1/2}$$

We will use function norms to measure the difference between the exact solution and the finite element solution. For example, we can define a particular measure of the error e as:

$$\|u - u_h\|_0 = \left(\int_{\Omega} (u - u_h)^2 d\Omega \right)^{1/2}$$

where u is the exact solution and u_h is the finite element solution.

How good is the FE solution?

Inevitably, the finite element solution is not exact. But how do we know that it bears any relation to the exact solution? If the mesh is refined, will the error reduce? We may say intuitively yes, but this is not very satisfactory. This is where a priori analysis helps. A priori error estimate:

$$\|u - u_h\|_s \leq Ch^\beta \|u\|_{k+1}$$

- s is the norm of interest,
- k is the polynomial order of the shape functions (complete polynomials),
- $\beta = \min(k + 1 - s, 2(k + 1 - m))$,
- m is the highest order derivative appearing in the weak form
- C is an unknown constant that does not depend on h or the exact solution u .

What is of key interest is the exponent β .

Error in the temperature/displacement field

For a steady heat conduction or elasticity problem, for the error in the temperature/displacement field:

- $s = 0$
- $m = 1$

Therefore, for $k \geq 1$,

$$\|u - u_h\|_s \leq Ch^{k+1} \|u\|_{k+1}$$

The solution converges with order $O(h^{k+1})$. Using quadratic shape functions, the error converges with order $O(h^3)$. So, if the element size is halved, the error will be reduced by a factor of eight.

Posterior error estimation and adaptivity

A *posterior error* estimation involves quantification of the error after a simulation has been performed.

If we have an a posterior error estimator for a particular problem, a simulation can be performed, the error estimated in various part of the domain and the mesh and finite element type can be adjusted in order to reduce the error to a prescribed value. This is known as *adaptivity*.

The two obvious strategies are known as:

- *h-adaptivity*: this involves modify the mesh.
- *p-adaptivity*: this involves changing the finite element type.

A *posterior* error estimation involves quantification of the error after a simulation has been performed.

If we have an a posterior error estimator for a particular problem, a simulation can be performed, the error estimated in various part of the domain and the mesh and finite element type can be adjusted in order to reduce the error to a prescribed value. This is known as adaptivity.

The two obvious strategies are known as:

- *h*-adaptivity: this involves modify the mesh.
- *p*-adaptivity: this involves changing the finite element type.

A *posterior* error estimation involves quantification of the error after a simulation has been performed. With an estimate of the error, a judgement can be made as to the reliability of the result. It is then also possible to adapt the finite element mesh to reduce the error. The goal of adaptivity is minimise the error for a given computational effort. It is a rich area of activity in applied and computational mathematics.

Error in the strain/flux field

For a steady heat conduction or elasticity problem, for the error in the strain/flux field we have:

- $s = 1$

- $m = 1$

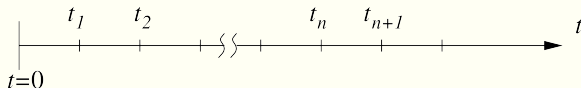
Therefore, for $k \geq 1$,

$$\|u - u_h\|_s \leq Ch^k \|u\|_{k+1}$$

For a given element, the order of convergence in the first derivative is one order less than for the unknown field itself. Hence, the stresses in a finite element simulation are usually 'less' accurate than the displacements.

Time-dependent problems

We wish to develop schemes to deal with semi-discrete finite element equations in a step-by-step fashion. Starting at time $t = 0$, we want to solve a problem up to time $t = t_{final}$. The time interval $(0, t_{final})$ will be 'chopped' into increments Δt :



Time increment: $\Delta t = t_{n+1} - t_n$

What is needed is a strategy to advance the solution from t_n to t_{n+1} .

There are many different techniques for this. We will consider some simple examples and the most commonly used methods.

Forward Euler

Given an equation of the form:

$$\dot{y} = f(y, t)$$

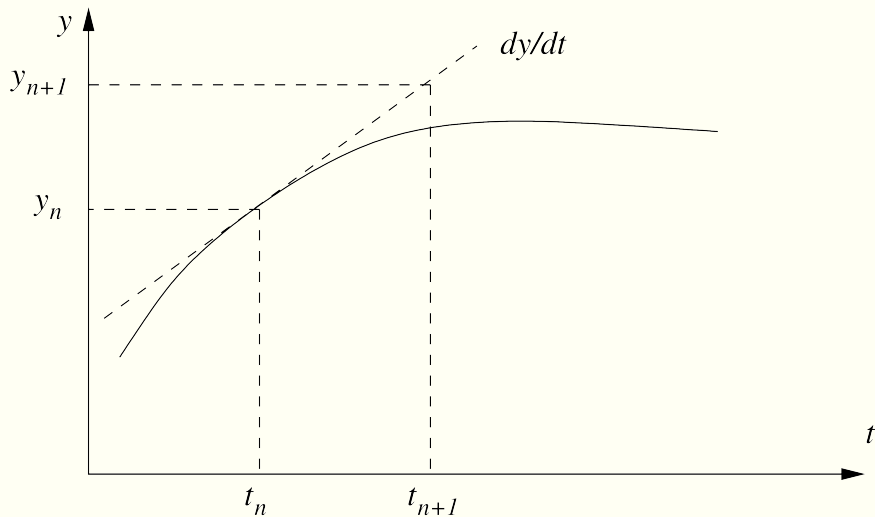
for which we know everything at time t_n (y_n and $f(y_n, t_n)$), we want to advance to time t_{n+1} and find y_{n+1} . The simplest approach is the forward Euler method. It involves:

$$y_{n+1} = y_n + \Delta t f(y_n, t_n)$$

Re-arranging,

$$\frac{y_{n+1} - y_n}{\Delta t} = f(y_n, t_n)$$

Forward Euler



Forward Euler

We can apply the forward Euler scheme to the nodal degrees of freedom \mathbf{a} by replacing y with \mathbf{a} and f with $\dot{\mathbf{a}}$,

$$\frac{\mathbf{a}_{n+1} - \mathbf{a}_n}{\Delta t} = \dot{\mathbf{a}}_n$$

We would like to eliminate the time derivative term $\dot{\mathbf{a}}$ from the semi-discrete heat equation, $\mathbf{M}\dot{\mathbf{a}} + \mathbf{K}\mathbf{a} = \mathbf{b}$, so that we are left with only \mathbf{a} . In the forward Euler expression, we have $\dot{\mathbf{a}}_n$, therefore we consider the heat equation at time step t_n ,

$$\mathbf{M}\dot{\mathbf{a}}_n + \mathbf{K}\mathbf{a}_n = \mathbf{b}_n$$

Inserting the forward Euler expression for $\dot{\mathbf{a}}$, we have

$$\mathbf{M}\frac{\mathbf{a}_{n+1} - \mathbf{a}_n}{\Delta t} + \mathbf{K}\mathbf{a}_n = \mathbf{b}_n$$

Taking all the known terms (those at t_n) to the right-hand side,

$$\mathbf{M}\mathbf{a}_{n+1} = \mathbf{M}\mathbf{a}_n + \Delta t(\mathbf{b}_n - \mathbf{K}\mathbf{a}_n)$$

Forward Euler

We have a system on linear equation that we can solve to find \mathbf{a}_{n+1} . Now the benefit of a lumped mass matrix can be seen. If \mathbf{M} is lumped, there is no need to solve a system of equations and finding \mathbf{a}_{n+1} is trivial.

System of equations for a lumped mass matrix:

$$\begin{bmatrix} M_{11} & 0 & 0 \\ 0 & M_{22} & 0 \\ 0 & 0 & M_{33} \end{bmatrix} \begin{bmatrix} a_{1,n+1} \\ a_{2,n+1} \\ a_{3,n+1} \end{bmatrix} = \begin{bmatrix} \hat{b}_{1,n} \\ \hat{b}_{2,n} \\ \hat{b}_{3,n} \end{bmatrix}$$

Solution:

$$a_{i,n+1} = \hat{b}_{i,n} / M_{ii}$$

This is the key to the efficient simulation of some large scale time-dependent finite element problems. We can start at $t = 0$ and advance in time. The process is simply repeated until the desired final time is reached.

CE394M: solvers - errors

- └ Time-dependent problems
 - └ Explicit schemes
 - └ Forward Euler

We have a system on linear equation that we can solve to find \mathbf{a}_{n+1} . Now the benefit of a lumped mass matrix can be seen. If \mathbf{M} is lumped, there is no need to solve a system of equations and finding \mathbf{a}_{n+1} is trivial.

System of equations for a lumped mass matrix:

$$\begin{bmatrix} M_{11} & 0 & 0 \\ 0 & M_{22} & 0 \\ 0 & 0 & M_{33} \end{bmatrix} \begin{bmatrix} a_{1,n+1} \\ a_{2,n+1} \\ a_{3,n+1} \end{bmatrix} = \begin{bmatrix} \tilde{b}_{1,n} \\ \tilde{b}_{2,n} \\ \tilde{b}_{3,n} \end{bmatrix}$$

Solution:

$$a_{i,n+1} = \tilde{b}_{i,n} / M_{ii}$$

This is the key to the efficient simulation of some large scale time-dependent finite element problems. We can start at $t = 0$ and advance in time. The process is simply repeated until the desired final time is reached.

The forward Euler scheme is an example of an explicit scheme since f depends only on data at t_n (or possibly earlier steps). The forward Euler scheme is very simple, but unfortunately not very useful for the heat equation. Firstly, it is not particularly accurate (only $O(\Delta t)$), although this is not the biggest issue. The problem is stability. If Δt is too large, the solution will grow rapidly and uncontrollably. The forward Euler method is notoriously unstable.

The time step at which a method becomes unstable is known as the *critical time step*. Conditionally stable methods are stable for: $\Delta t \leq \Delta t_{\text{crit}}$

- The critical time step often scales with eigenvalues of the finite element system, which in turn depend on the element type and size.
- As the element size is reduced, the critical time step reduces.
- The heat equation is a so-called stiff equation. This means that the restriction on the time step for explicit schemes is particularly severe. It scales according to $\Delta t \propto h^2$

where h is a measure of the size of an element. This is a very strong restriction, and means that explicit schemes are not usually practical for the heat equation.

Backward Euler

$$y_{n+1} = y_n + \Delta t f(y_{n+1}, t_{n+1})$$

Note now that f is evaluated at y_{n+1} and not at y_n . However, we do not know y_{n+1} ; it is what we want to find.

$$\begin{aligned}\mathbf{a}_{n+1} &= \mathbf{a}_n + \Delta t \dot{\mathbf{a}}_{n+1} \\ \frac{\mathbf{a}_{n+1} - \mathbf{a}_n}{\Delta t} &= \dot{\mathbf{a}}_{n+1}\end{aligned}$$

Since the time derivative term is evaluated at t_{n+1} , we pose the semi-discrete heat equation at t_{n+1} ,

$$\mathbf{M} \dot{\mathbf{a}}_{n+1} + \mathbf{K} \mathbf{a}_{n+1} = \mathbf{b}_{n+1}$$

Inserting the backward Euler expression for $\dot{\mathbf{a}}$, we have

$$\mathbf{M} \frac{\mathbf{a}_{n+1} - \mathbf{a}_n}{\Delta t} + \mathbf{K} \mathbf{a}_{n+1} = \mathbf{b}_{n+1}$$

Collecting the unknown terms on the left-hand side and the known terms on the right- hand side:

$$(\mathbf{M} + \Delta t \mathbf{K}) \mathbf{a}_{n+1} = \Delta t \mathbf{b}_{n+1} + \mathbf{M} \mathbf{a}_n$$

- Now, even if \mathbf{M} is a diagonal matrix, we still need to solve a system of equations to find \mathbf{a}_{n+1} because \mathbf{K} is not diagonal.
- A system must be solved because the backward Euler method is an *implicit method*.
- In stepping from t_n to t_{n+1} , it relies on data at t_{n+1} which is what characterises implicit methods.
- The consequence is that a system of equations must be solved.

CE394M: solvers - errors

- └ Time-dependent problems
 - └ Implicit method
 - └ Backward Euler

Collecting the unknown terms on the left-hand side and the known terms on the right-hand side:

$$(\mathbf{M} + \Delta t \mathbf{K}) \mathbf{a}_{n+1} = \Delta t \mathbf{b}_{n+1} + \mathbf{M} \mathbf{a}_n$$

- ◆ Now, even if \mathbf{M} is a diagonal matrix, we still need to solve a system of equations to find \mathbf{a}_{n+1} because \mathbf{K} is not diagonal.
- ◆ A system must be solved because the backward Euler method is an implicit method.
- ◆ In stepping from t_n to t_{n+1} , it relies on data at t_{n+1} which is what characterizes implicit methods.
- ◆ The consequence is that a system of equations must be solved.

The backward Euler method requires far more computational effort at each time step compared to the forward Euler method, yet it has the same order of accuracy ($O(\Delta t)$). However it is unconditionally stable. That is, there is no critical time step and the solution will remain bounded for any Δt . While each step is more expensive relative to an explicit method, it is possible that far fewer total steps will be required.

Summary of time-stepping schemes

- Explicit schemes are characterised by conditional stability.
- Explicit schemes can sometimes be very efficient as there is no need to solve systems of equations.
- Implicit schemes can be unconditionally stable (although not all implicit schemes are unconditionally stable).
- Implicit schemes require the solution of a system of equations, which can be expensive.
- Whether a scheme is explicit or implicit does not infer anything regarding accuracy.