

## 7. Galerkin Methods.

### 7.1 Introduction.

Until now, finite-difference methods for solving partial differential equations were applied. These methods specify the dependent variables at certain gridpoints in space and time and the derivatives in the equations are evaluated using Taylor Series expansions. The definitions of convergence, accumulated error, ... are based on comparing the solution  $U_j^n = U(j \Delta x, n \Delta t)$  to the continuous solution  $u(x, t)$  at grid point locations. The Galerkin procedure represents the dependent variables with a sum of functions that have a prescribed spatial structure. The coefficient associated with each function is normally a function of time. This procedure transforms a partial differential equation into a set of ordinary differential equations for the coefficients which are usually solved with finite differences in time. The two most useful Galerkin methods are the spectral method and the finite element method.

We now look at grid point values  $U_j^n$  as being representative of grid-box averages of  $u(x, t)$ . Thus, in the case of one spatial dimension, we now compare

$$U_j^n = \frac{1}{\Delta x} \int_{(j-1/2)\Delta x}^{(j+1/2)\Delta x} U_j^n dx \quad \text{to} \quad \frac{1}{\Delta x} \int_{(j-1/2)\Delta x}^{(j+1/2)\Delta x} u(x, t) dx$$

rather than  $U_j^n$  to  $u(x \Delta x, n \Delta t)$ .

We can reformulate by defining

$$(1) \quad \phi_j(x) = \begin{cases} 1/\Delta x & \text{for } (j-1/2)\Delta x \leq x < (j+1/2)\Delta x \\ 0 & \text{elsewhere} \end{cases}$$

and the  $U_j^n$  values <sup>as a vehicle for defining</sup> a piecewise constant approximation  $U(x, n\Delta t)$  such that

$$(2) \quad U(x, n\Delta t) = \sum_{j=-\infty}^{\infty} U_j^n \phi_j(x)$$

A formal series expansion for the purpose of approximating  $u(x, t)$  can be carried out with an infinite variety of functions  $\phi_j$ . The ones from (1) are representative of the nodal values associated with the standard finite-difference equations.

Example: Advection equation.  $\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$

By using (2), the advection equation can be rewritten

$$(3) \quad \sum_j \frac{\partial U_j^n}{\partial t} \phi_j = -c \sum_j U_j^n \frac{\partial \phi_j}{\partial x}$$

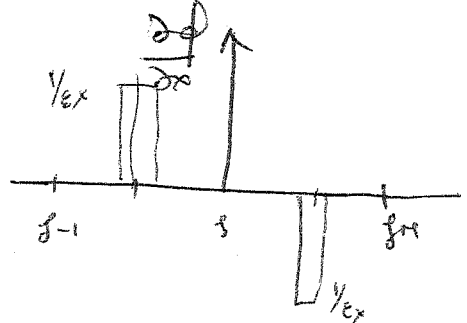
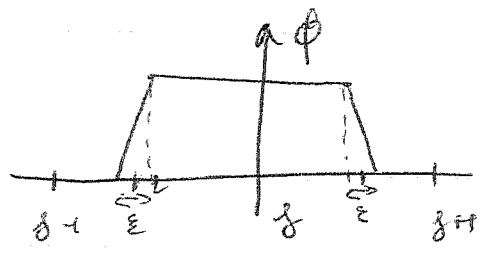
The basis functions (1) are orthogonal  $\Rightarrow$

$$\int_{-\infty}^{\infty} \phi_i(x) \phi_j(x) dx = \begin{cases} (\Delta x) & \text{if } i=j \\ 0 & \text{elsewhere} \end{cases}$$

We can obtain an equation for  $\frac{\partial U_i^n}{\partial t}$  by multiplying by  $\phi_k$  and integrating over  $x$

$$(4) \quad \sum_j \frac{\partial U_j^n}{\partial t} \int_{-\infty}^{\infty} \phi_j \phi_k dx = -c \sum_j U_j^n \int_{-\infty}^{\infty} \frac{\partial \phi_j}{\partial x} \phi_k dx$$

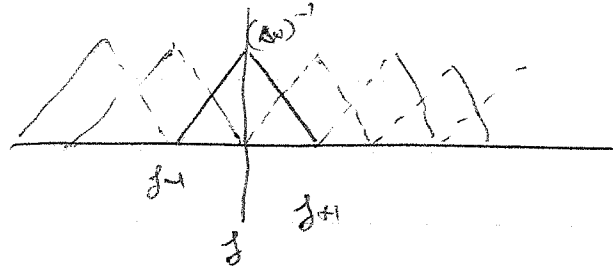
In order to compute the RBS integral, we consider the step function to be a trapezoid



$$\Rightarrow \int_{-\infty}^{\infty} \frac{\partial \phi_j}{\partial x} \phi_k dx = \begin{cases} -\frac{1}{2\Delta x} & j = k-1 \\ +\frac{1}{2\Delta x} & j = k+1 \\ 0 & \text{elsewhere} \end{cases} \quad \left( \begin{array}{l} \text{independent} \\ \text{of the slope} \end{array} \right) \quad (\text{valid for } \epsilon \ll \Delta x)$$

(4) then becomes  $\frac{\partial U_k^n}{\partial t} = -c \frac{U_{k+1}^n - U_{k-1}^n}{2\Delta x}$  centered in space

\* If we use a different  $\phi$  such as a piecewise linear representation



They are not orthogonal  $\int_{-\infty}^{\infty} \phi_j \phi_k dx = \begin{cases} \frac{2}{3\Delta x} & j = k \\ \frac{1}{6\Delta x} & |j-k|=1 \\ 0 & \text{elsewhere} \end{cases}$

and result in

(5)  $\frac{1}{6} \left( \frac{\partial U_{j-1}^n}{\partial t} + 4 \frac{\partial U_j^n}{\partial t} + \frac{\partial U_{j+1}^n}{\partial t} \right) = -c \frac{U_{k+1}^n - U_{k-1}^n}{2\Delta x}$

\* We can generalize these ideas into a formal definition of the Galerkin method

Given a differential equation  $L(u) = f(x)$  where  $L$  is a differential operator,  $u$  the dependent variable and  $f(x)$  a specified forcing function.

in the domain  $R$  ( $x$  may be multidimensional).  
The Galerkin approximation is defined by

$$(6) \quad U(x, t) = \sum_{j=1}^N A_j(t) \phi_j(x) \quad \varepsilon$$

where the coefficients  $A_j(t)$  are determined by requiring that the error

$$(7) \quad e_N = \langle U(x, t) - f(x) = \langle \left( \sum_{j=1}^N A_j(t) \phi_j(x) \right) - f(x) \rangle$$

be orthogonal to each basis function.

$$(8) \quad \int_R e_N \phi_j(x) dx = 0 \quad j=1, \dots, N$$

$\underbrace{\hspace{10em}}_{\text{domain}}$

The final form is

$$(9) \quad \int_R \phi_k \left( \sum_{j=1}^N A_j(t) \phi_j(x) \right) dx - \int_R \phi_k f(x) dx = 0$$

$k=1, \dots, N.$

This reduces to the problem of  $N$  algebraic equations that relate the unknown coefficients  $A_j(t)$  to the "transform" of the forcing function. They are normally solved by finite difference in time.

There are various ways to interpret (8).

- (1) The residual error is orthogonal to orthogonal to  $\phi_j$ , i.e. the error should have no components in the space spanned by the  $\phi_j$ .

(2) The coefficients  $A_j$  should be chosen to minimize the integral  $\int_R e^2(x,t) dx$

Straight forward when  $L$  is a linear operator  
In more complicated cases, not necessarily valid.

(3)  $L(u) = f(x)$  is approximated by  $L(U) = f(x)$   
as in the Introduction.

Schemes employing basis functions defined in terms of periodic functions are referred to as "spectral". The ones using more "localized" basis functions are "finite elements" schemes.

### 7.2 Energy conservation

If we consider the simplified equation

(10) 
$$\frac{\partial u}{\partial t} + L(u) = 0$$

then the Galerkin form is

(11) 
$$\sum_{j=1}^N \frac{\partial A_j}{\partial t} \int_R \phi_k \phi_j dx + \int_R \phi_k L \left( \sum_{j=1}^N A_j \phi_j \right) dx = 0$$
  
$$k=1, \dots, N$$

This process gives  $N$  coupled ordinary differential equations in the coefficients  $A_j(t)$ . This can be solved by introducing finite differences in time.

We already discussed the importance of energy conserving schemes. The Galerkin method leads naturally to energy conservation in equations with quadratic energy invariants.

For an energy conserving system

$$(12) \quad \int_R \frac{\partial (u^2/2)}{\partial t} = - \int_R u \mathcal{L}(u) dx$$

, the operator  $\mathcal{L}$  must satisfy the condition

$$\int_R \psi \mathcal{L}(\psi) dx = 0 \quad \text{where } \psi \text{ is any}$$

reasonable function that satisfies the boundary conditions.

Then (12) satisfies

$$(13) \quad \frac{d}{dt} \int_R u^2/2 dx = 0$$

which shows the energy conservation for the exact equation. We need to demonstrate that it holds for the finite sum.

We multiply the  $k^{\text{th}}$  equation (11) by  $A_k$  and sum from  $k=1$  to  $N$ .

$$(14) \quad \int_R \left( \sum_{k=1}^N A_k \phi_k \right) \frac{\partial}{\partial t} \left( \sum_{s=1}^N A_s \phi_s \right) dx = \\ - \int_R \left( \sum_{k=1}^N A_k \phi_k \right) \mathcal{L} \left( \sum_{s=1}^N A_s \phi_s \right) dx$$

The integral on the right side vanishes if

$$\text{we set } \psi = \sum_{s=1}^N A_s \phi_s = \sum_{k=1}^N A_k \phi_k$$

and (14) can be rewritten

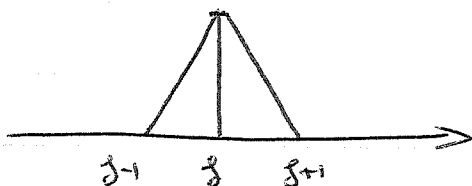
$$\frac{1}{2} \int_R \frac{\partial}{\partial t} \left( \sum_{k=1}^N A_k \phi_k \right)^2 dx = 0$$

Energy conservation  
for the Galerkin  
Approximation

### 7.3 The Advection equation with Finite Elements

$$\frac{\partial u}{\partial t} + c \frac{\partial u}{\partial x} = 0$$

We use regular piecewise linear elements such as



$$\phi_j(x) = \begin{cases} 0 & x < (j-1)\Delta x \\ (x - (j-1)\Delta x)/\Delta x & (j-1)\Delta x < x < j\Delta x \\ ((j+1)\Delta x - x)/\Delta x & j\Delta x < x < (j+1)\Delta x \\ 0 & x > (j+1)\Delta x \end{cases}$$

The Galerkin equation is obtained by setting  $L = c \frac{\partial}{\partial x}$ ,  $f(x) = 0$

$$(5) \quad \sum_{j=1}^N \frac{\partial A_j}{\partial t} \int_R \phi_k \phi_j dx + c \sum_{j=1}^N A_j \int_R \phi_k \frac{\partial \phi_j}{\partial x} dx = 0$$

$k = 1, \dots, N$

The resulting equation is the one derived in the introduction

$$(16) \quad \frac{1}{6} \left( \frac{\partial A_{k+1}}{\partial t} + 4 \frac{\partial A_k}{\partial t} + \frac{\partial A_{k-1}}{\partial t} \right) = -c \frac{A_{k+1} - A_{k-1}}{2\Delta x}$$

The advection term is the same as if obtained from centered differencing, but the time derivatives appears as a weighted average over three points. This greatly increases the accuracy of the solution.

We now apply the leap-frog time differencing scheme

$$(17) \quad \frac{1}{12\Delta t} \left( A_{k+1}^{n+1} - A_{k+1}^{n-1} + 4(A_k^{n+1} - A_k^{n-1}) + A_{k-1}^{n+1} - A_{k-1}^{n-1} \right) = -c \frac{A_{k+1}^n - A_{k-1}^n}{2\Delta x}$$

The stability and phase error of this scheme can be investigated by substituting  $A_k^n = A e^{i(\mu\Delta x k + \alpha n \Delta t)}$  (Variation of the Fourier Transform).

Substitution into (17) leads to

$$(18) \quad \sin \alpha \Delta t = -\frac{c\Delta t}{\Delta x} \left( \frac{3 \sin \mu \Delta x}{k} + \cos \mu \Delta x \right)$$

(For the leap frog  $\sin \alpha \Delta t = -\frac{c\Delta t}{\Delta x}$ )

The solution is stable (neutral solution with no damping or amplification) if  $\alpha$  is real or  $|\sin \alpha \Delta t| \leq 1$ . To ensure stability for all wavelengths, it is necessary to find the maximum magnitude of the RHS of (18)

Maximum when  $\mu \Delta x = 120^\circ \Rightarrow$

$$|c \Delta t / \Delta x| \leq 1/\sqrt{3}$$

which is more restrictive than the leapfrog FD scheme. However it gives even better phase speed than the fourth-order leapfrog scheme

Finite elements are an interesting alternative to classic FD methods. They offer a high level of flexibility offered for the use of grids of



variable size, shape and flexibility and are attractive despite a higher cost in computer time. They are popular in the engineering field and control domains. For a review, see Provost (1985) in O'Brien

7. The Spectral and Pseudo-Spectral method applied to the <sup>non-linear</sup> advection equation (Burgers)

Before the advent of the fast Fourier Transform (FFT) spectral methods played only a minor role in fluid dynamics because they were far less economical than grid points methods.

We want to solve  $u_t(u) = \frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x}$

Let us use a cyclic boundary condition on the domain  $-1 \leq x \leq 1$

We use for the basis functions, the trigonometric

(19)  $\phi_j(x) = e^{i\pi j x} \quad (j = -J, \dots, J)$

The Galerkin approximation is then

(20)  $\langle \phi_m, \sum_j A'_j \phi_j \rangle + \langle \phi_m, (\sum_j A_j \phi_j) (\sum_k A_k \phi'_k) \rangle = 0$   
 $m = -J, \dots, J$

and  $\langle \phi_i, \phi_j \rangle = \int_R \phi_i \bar{\phi}_j$  (conjugate (imaginary numbers))

The general form of this set of equations is

(21)  $\sum_j a_{mj} A'_j + \sum_j \sum_k b_{mj k} A_j A_k = 0 \quad m = -J \rightarrow J$

To advance the solution by one step, one needs  $(2\sigma+1)^3$  multiplications for the "interaction" term (RHS of (21)) and an inversion of the matrix of capacitances  $a_{mj}$ . The basis functions are orthogonal  $\Rightarrow$

$$\langle \phi_m, \phi_j \rangle = \begin{cases} 2 & m=j \\ 0 & \text{otherwise} \end{cases}$$

$$\langle \phi_m, \phi_j \phi'_k \rangle = \begin{cases} 2i\pi k & m=jk \\ 0 & \text{otherwise} \end{cases}$$

(21) then reduces to

$$(22) \quad A'_m + \sum_{j,k=m} i\pi k A_j A_k = 0$$

Thus, in practice, the inversion problem does not arise and the operation count for the interaction terms is only  $\sigma^2$ . The operation count for grid point methods on the other is proportional to the number of grid points  $(2\sigma+1)$ . The difference is significant and acted as a major deterrent in the past.

The FFT can improve the speed of the spectral method. There is not apparent need to perform Fourier transforms except at  $t=0$  when the initial conditions have to be transformed from physical to phase space and vice-versa at the end. The speed of these transforms does not affect the overall efficiency of the method.

The so-called "pseudo-spectral" method produces similar results to the spectral method by transforming the variables back and forth between grid point and phase space every few steps.

The transform (Pseudo-spectral) method sees the series at certain spatial grid points and these fields are multiplied together at each point to form the non linear terms. Then these terms are transformed back to spectral space. The usefulness of this method is enhanced by the existence of efficient transform methods such as the FFTs. This method is essentially a grid point method which uses spectral decomposition techniques to eliminate the problem of FD in space, namely the phase retardation of short waves by capturing the spatial derivatives by differentiating the individual Fourier components.

We define the grid points by  $x_j = j\Delta x$  for  $-\sigma < j < \sigma$ . With  $\phi_k(x_j) = e^{i\pi k x_j}$  as basis functions we have

$$(23) \quad U_j(t) = U(x_j, t) = \sum_{k=-\sigma}^{\sigma-1} A_k(t) e^{i\pi k x_j}$$

Note that the number of Fourier components matches the number of grid points (Cyclic). Also, if  $k = \sigma$ , then we have two identical basis functions.

The orthogonality relation is then

$$(24) \quad \langle \phi_p, \phi_q \rangle = \frac{1}{\sigma} \sum_{j=-\sigma}^{\sigma-1} e^{i\pi(p-q)x_j} = \begin{cases} 2 & p=q \text{ mod } \sigma \\ 0 & \text{otherwise} \end{cases} \quad (25)$$

The inverse transform of (23) is then

$$(25) \quad A_k(t) = \frac{1}{2\sigma} \sum_{j=-\sigma}^{\sigma-1} U_j(t) e^{-i\pi k x_j}$$

Starting with grid points values  $U_j$ , the

The coefficients  $A_k$  can be computed from (25). The derivative  $\partial U_3 / \partial x$  is obtained by a second transform

$$(26) \quad \frac{\partial U_3(t)}{\partial x} = \sum_{k=-J}^{J-1} i\pi k A_k(t) e^{i\pi k x_j}$$

Time integration is normally done by finite-differencing. An alternative method is to perform the time integration in phase space, but then transform the dependent variables back and forth to grid point space for the evaluation of the nonlinear terms.

The final equation which approximates the advection equation is then

$$(27) \quad \underbrace{\frac{\partial U_3(t)}{\partial t}}_C + \underbrace{U_3(t)}_{\text{Basis}} \sum_{k=-J}^{J-1} i\pi k A_k(t) e^{i\pi k x_j} = 0$$

Let's now compare the two methods (Spectral and Pseudo Spectral).

We first transform (27) (PS) into an equation for the  $A_k$

$$(28) \quad (25) A'_k(t) + \sum_P \sum_Q \sum_I i\pi q A_P(t) A_Q(t) e^{i\pi(P+Q-k)x_j} = 0$$

Using the orthogonality of the basis functions, we can rewrite (28) as

$$(29) \quad A_k^I(t) + \sum_{p+q=k} i\pi q A_p A_q + \sum_{p+q=k+2J} i\pi q A_p A_q + \sum_{p+q=k-2J} i\pi q A_p A_q = 0 \quad (\text{since the orthogonality is modulo } 2J)$$

The last two sums show a significant difference between (P) and (PS). They are referred to as "aliased" terms brought by the finite sampling interval of the discrete Fourier decomposition. This method is clearly much faster since the number of operations is  $(2J+1) \log_2(2J+1)$  versus  $(2J+1)^2$  for the unaccelerated method (P).

There are two basic techniques for removing the aliasing error introduced in (29)

(1) Aliasing removal by Padding or Truncation. The key is to use a discrete transform with  $M$  rather than  $N$  points where  $M \gg \frac{3N}{2}$ .

(2) Aliasing removal by phase shifts

Both methods can be extended to two and three dimensions. For a complete description of spectral methods

"Spectral Methods in Fluid Dynamics"

by Canuto

Hussaini

Quarteroni

Zang

Springer-Verlag

1988