

Math 6630: Numerical Solutions of Partial Differential Equations
Fourier spectral methods, I
See Canuto et al. 2011, Chapter 3,
Shen, Tang, and Wang 2011, Chapter 2

Akil Narayan¹

¹Department of Mathematics, and Scientific Computing and Imaging (SCI) Institute
University of Utah

March 22, 2023



Weighted residual methods

We can now discuss some basic approaches for constructing a particular class of weighted residual methods: spectral methods.

Recall: with \mathcal{L} a (time-independent, stationary) linear operator and f a given function, we seek to compute the solution u to

$$\mathcal{L}(u) = f, \quad \mathcal{R}(u) := \mathcal{L}(u) - f. \quad + \text{BC's}$$

Strong notions of convergence are too restrictive, so we employ weak notions:

$$\langle \mathcal{R}(u), v \rangle = 0, \quad \forall v \in V,$$

where $U \ni u$ is a chosen trial space, and $V \ni v$ is a chosen test space.

These are generic [weighted residual methods](#).

The inner product $\langle \cdot, \cdot \rangle$ corresponds to a particular Hilbert space (for us, almost always L^2).

Weighted residual methods

We can now discuss some basic approaches for constructing a particular class of weighted residual methods: spectral methods.

Recall: with \mathcal{L} a (time-independent, stationary) linear operator and f a given function, we seek to compute the solution u to

$$\mathcal{L}(u) = f, \quad \mathcal{R}(u) := \mathcal{L}(u) - f. \quad + \text{BC's}$$

Strong notions of convergence are too restrictive, so we employ weak notions:

$$\langle \mathcal{R}(u), v \rangle = 0, \quad \forall v \in V,$$

where $U \ni u$ is a chosen trial space, and $V \ni v$ is a chosen test space.

These are generic [weighted residual methods](#).

The inner product $\langle \cdot, \cdot \rangle$ corresponds to a particular Hilbert space (for us, almost always L^2).

We'll discuss methods where we choose U and/or V as a space of Fourier functions.

This will be a particular instance of a [spectral method](#).

The Fourier-Galerkin method

The Fourier-Galerkin method:

- Is a Galerkin procedure: $U = V$
- Uses (global) Fourier functions: $V = \text{span}\{e^{ikx} \mid k \in \mathbb{Z}\}$ (assuming one spatial dimension on $[0, 2\pi]$)

The Fourier-Galerkin method

The Fourier-Galerkin method:

- Is a Galerkin procedure: $U = V$
- Uses (global) Fourier functions: $V = \text{span}\{e^{ikx} \mid k \in \mathbb{Z}\}$ (assuming one spatial dimension on $[0, 2\pi]$)

Of course, this choice only makes sense if the boundary conditions are periodic.

Recall that a proper strategy to address this problem is to identify the appropriate bilinear form a :

$$\mathcal{L}(u) = f \xrightarrow{\text{IbP}} a(u, v) = \langle f, v \rangle.$$

The Fourier-Galerkin method

The Fourier-Galerkin method:

- Is a Galerkin procedure: $U = V$
- Uses (global) Fourier functions: $V = \text{span}\{e^{ikx} \mid k \in \mathbb{Z}\}$ (assuming one spatial dimension on $[0, 2\pi]$)

Of course, this choice only makes sense if the boundary conditions are periodic.

Recall that a proper strategy to address this problem is to identify the appropriate bilinear form a :

$$\mathcal{L}(u) = f \xrightarrow{\text{IbP}} a(u, v) = \langle f, v \rangle.$$

One can then choose finite-dimensional spaces U_N and V_N for discretization. I.e.,:

Find $u_N \in U_N$ satisfying $a(u_N, v) = \langle f, v \rangle$, for all $v \in V$. U_N

An elliptic equation

Let's construct a Fourier-Galerkin scheme for the following problem:

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's}$$

Since this is a Galerkin scheme, $U = V$, and so we need only identify V . With periodic boundary conditions, it's sensible to choose:

$$u, v \in \text{span}\{\phi_k(x)\}_{k \in \mathbb{Z}}, \quad \phi_k(x) = \frac{1}{\sqrt{2\pi}} e^{ikx}.$$

Using an L^2 inner product on $[0, 2\pi]$, then

$$\langle -u'' + u, v \rangle \xrightarrow{\text{IbP}} \langle u', v' \rangle + \langle u, v \rangle.$$

An elliptic equation

Let's construct a Fourier-Galerkin scheme for the following problem:

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's}$$

Since this is a Galerkin scheme, $U = V$, and so we need only identify V . With periodic boundary conditions, it's sensible to choose:

$$u, v \in \text{span}\{\phi_k(x)\}_{k \in \mathbb{Z}}, \quad \phi_k(x) = \frac{1}{\sqrt{2\pi}} e^{ikx}.$$

Using an L^2 inner product on $[0, 2\pi]$, then

$$\langle -u'' + u, v \rangle \xrightarrow{\text{IbP}} \langle u', v' \rangle + \langle u, v \rangle.$$

This immediately suggests an appropriate sesquilinear form:

$$a(u, v) := \langle u', v' \rangle + \langle u, v \rangle,$$

$|a(v, v)| \geq \|v\|_V^2$
 $|a(u, v)| \leq \sqrt{2} \|u\|_V \|v\|_V$

which we know is continuous and coercive with respect to H_p^1 .

Thus, we choose $V = H_p^1$, and the dual of V with respect to L^2 is $V^* = H_p^{-1}$.

Thus (V, L^2, V^*) is our Gelfand triple, and for any $f \in V^*$,

$$\text{Find } u \in V \text{ satisfying } a(u, v) = \langle f, v \rangle \text{ for all } v \in V$$

is well-posed by Lax-Milgram.

A Fourier-Galerkin scheme

Find $u \in V$ satisfying $a(u, v) = \langle f, v \rangle$ for all $v \in V$

We computed last time:

$$\sup_{u, v \in V} \frac{|a(u, v)|}{\|u\|_V \|v\|_V} \leq \sqrt{2} =: C, \quad \inf_{v \in V} \frac{|a(v, v)|}{\|v\|_V^2} = 1 =: c.$$

which will be useful when computing the error in our scheme.

A Fourier-Galerkin scheme

Find $u \in V$ satisfying $a(u, v) = \langle f, v \rangle$ for all $v \in V$

We computed last time:

$$\sup_{u, v \in V} \frac{|a(u, v)|}{\|u\|_V \|v\|_V} \leq \sqrt{2} =: C, \quad \inf_{v \in V} \frac{|a(v, v)|}{\|v\|_V^2} = 1 =: c.$$

which will be useful when computing the error in our scheme.

To turn all this into a scheme, we need to discretize. Let's choose,

$$V_N = \text{span} \{ \phi_k, \mid |k| \leq N \} \subset V.$$

$$\dim V_N = 2N + 1$$

Then our discrete scheme is

Find $u_N \in V_N$ satisfying $a(u_N, v) = \langle f, v \rangle$ for all $v \in V_N$

A Fourier-Galerkin scheme

Find $u \in V$ satisfying $a(u, v) = \langle f, v \rangle$ for all $v \in V$

We computed last time:

$$\sup_{u, v \in V} \frac{|a(u, v)|}{\|u\|_V \|v\|_V} \leq \sqrt{2} =: C, \quad \inf_{v \in V} \frac{|a(v, v)|}{\|v\|_V^2} = 1 =: c.$$

which will be useful when computing the error in our scheme.

To turn all this into a scheme, we need to discretize. Let's choose,

$$V_N = \text{span} \{ \phi_k, \mid |k| \leq N \} \subset V.$$

Then our discrete scheme is

Find $u_N \in V_N$ satisfying $a(u_N, v) = \langle f, v \rangle$ for all $v \in V_N$

The above problem is equivalent to making the ansatz,

$$u_N(x) = \sum_{|k| \leq N} \hat{u}_k \phi_k(x),$$

and using $v \leftarrow \phi_k$ for every k satisfying $|k| \leq N$.

Scheme details

Find $u_N \in V_N$ satisfying $a(u_N, v) = \langle f, v \rangle$ for all $v \in V_N$

$$u_N(x) = \sum_{|k| \leq N} \hat{u}_k \phi_k(x),$$

With $v = \phi_k$ for a fixed k , the weak form reads,

$$\left\langle \sum_{|\ell| \leq N} (i\ell) \hat{u}_\ell \phi_\ell, (ik) \phi_k \right\rangle + \langle u_N, \phi_k \rangle = \langle f, \phi_k \rangle.$$

Scheme details

Find $u_N \in V_N$ satisfying $a(u_N, v) = \langle f, v \rangle$ for all $v \in V_N$

$$u_N(x) = \sum_{|k| \leq N} \hat{u}_k \phi_k(x),$$

With $v = \phi_k$ for a fixed k , the weak form reads,

$$\left\langle \sum_{|\ell| \leq N} (i\ell) \hat{u}_\ell \phi_\ell, (ik) \phi_k \right\rangle + \langle u_N, \phi_k \rangle = \langle f, \phi_k \rangle.$$

Then defining,

$$\hat{f}_k = \langle f, \phi_k \rangle, \quad d_k = k^2 + 1,$$

our scheme is,

$$\widehat{\mathbf{D}}_2 \hat{\mathbf{u}} = \hat{\mathbf{f}}, \quad \widehat{\mathbf{D}}_2 = \text{diag} \left((-N)^2 + 1, (-N + 1)^2 + 1, \dots, (N - 1)^2 + 1, (N^2) + 1 \right),$$

and so $\hat{\mathbf{u}} = \widehat{\mathbf{D}}^{-1} \hat{\mathbf{f}}$ is the solution, prescribing u_N .

$\widehat{\mathbf{D}}_2$ is sometimes called a “modal” differentiation matrix.

In practice, we cannot compute \hat{f}_k exactly, so these coefficients are approximated with quadrature/interpolation (injecting aliasing error into the right-hand side).

Error estimates

One of the significant advantages of all this is that we have already laid all the groundwork we need to compute a direct error estimate.

First, we have that,

$$\|u - u_N\|_V \leq \frac{C}{c} \inf_{v \in V_N} \|u - v\|_V, \quad \frac{C}{c} = \sqrt{2},$$

due to [Céa's Lemma](#), which in turn relies on [Lax-Milgram](#).

Error estimates

One of the significant advantages of all this is that we have already laid all the groundwork we need to compute a direct error estimate.

First, we have that,

$$\|u - u_N\|_V \leq \frac{C}{c} \inf_{v \in V_N} \|u - v\|_V, \quad \frac{C}{c} = \sqrt{2},$$

due to [Céa's Lemma](#), which in turn relies on [Lax-Milgram](#).

We also have,

$$\inf_{v \in v_N} \|u - v\|_V \leq \|u - P_N u\|_V,$$

where P_N is the L^2 -orthogonal projection onto $\cancel{W} = \mathbb{H}_0^1 = V_N$.

Error estimates

One of the significant advantages of all this is that we have already laid all the groundwork we need to compute a direct error estimate.

First, we have that,

$$\|u - u_N\|_V \leq \frac{C}{c} \inf_{v \in V_N} \|u - v\|_V, \quad \frac{C}{c} = \sqrt{2},$$

due to [Céa's Lemma](#), which in turn relies on [Lax-Milgram](#).

We also have,

$$\inf_{v \in v_N} \|u - v\|_V \leq \|u - P_N u\|_V,$$

where P_N is the L^2 -orthogonal projection onto $V = H_p^1$.

Finally, we have our basic approximation estimate for Fourier series:

$$\|u - P_N u\|_{H^1} \leq N^{1-s}. \quad u \in H_p^s \Rightarrow \|u - P_N u\|_{H_p^r} \leq N^{r-s}$$

Error estimates

One of the significant advantages of all this is that we have already laid all the groundwork we need to compute a direct error estimate.

First, we have that,

$$\|u - u_N\|_V \leq \frac{C}{c} \inf_{v \in V_N} \|u - v\|_V, \quad \frac{C}{c} = \sqrt{2},$$

due to [Céa's Lemma](#), which in turn relies on [Lax-Milgram](#).

We also have,

$$\inf_{v \in v_N} \|u - v\|_V \leq \|u - P_N u\|_V,$$

where P_N is the L^2 -orthogonal projection onto $V = H_p^1$.

Finally, we have our basic approximation estimate for Fourier series:

$$\|u - P_N u\|_{H^1} \leq N^{1-s}.$$

One final piece of information: while Lax-Milgram establishes $f \in H_p^{-1} \Rightarrow u \in H_p^1$, one can also show that $f \in H_p^r \Rightarrow u \in H_p^{r+2}$.

Error estimates

One of the significant advantages of all this is that we have already laid all the groundwork we need to compute a direct error estimate.

First, we have that,

$$\|u - u_N\|_V \leq \frac{C}{c} \inf_{v \in V_N} \|u - v\|_V, \quad \frac{C}{c} = \sqrt{2},$$

due to [Céa's Lemma](#), which in turn relies on [Lax-Milgram](#).

We also have,

$$\inf_{v \in v_N} \|u - v\|_V \leq \|u - P_N u\|_V,$$

where P_N is the L^2 -orthogonal projection onto $V = H_p^1$.

Finally, we have our basic approximation estimate for Fourier series:

$$\|u - P_N u\|_{H^1} \leq N^{1-s}.$$

One final piece of information: while Lax-Milgram establishes $f \in H_p^{-1} \Rightarrow u \in H_p^1$, one can also show that $f \in H_p^r \Rightarrow u \in H_p^{r+2}$.

Thus, overall, we have that for any $r > -1$,

$$f \in H_p^r \implies u \in H_p^{r+2} \implies \|u - u_N\|_{H^1} \leq \sqrt{2} N^{1-(r+2)} \lesssim N^{-(r+1)}.$$

This convergence rate is optimal (since it is the optimal rate of approximation).

A “strong” form Fourier-Galerkin method

Sometimes in practice a “strong” version of the Fourier-Galerkin method will be implemented. This approach directly imposes the Galerkin condition on the strong PDE form:

$$\langle -u''_N + u, v \rangle = \langle f, v \rangle, \quad v \in V.$$

A “strong” form Fourier-Galerkin method

Sometimes in practice a “strong” version of the Fourier-Galerkin method will be implemented. This approach directly imposes the Galerkin condition on the strong PDE form:

$$\langle -u_N'' + u, v \rangle = \langle f, v \rangle, \quad v \in V.$$

Note that this is not unreasonable since if u_N is smooth enough, we have,

$$\langle -u_N'', v \rangle = \langle u_N', v' \rangle,$$

which implies that the strong form above is exactly equivalent to our first weak formulation that employed bilinear forms.

A “strong” form Fourier-Galerkin method

Sometimes in practice a “strong” version of the Fourier-Galerkin method will be implemented. This approach directly imposes the Galerkin condition on the strong PDE form:

$$\langle -u_N'' + u, v \rangle = \langle f, v \rangle, \quad v \in V.$$

Note that this is not unreasonable since if u_N is smooth enough, we have,

$$\langle -u_N'', v \rangle = \langle u_N', v' \rangle,$$

which implies that the strong form above is exactly equivalent to our first weak formulation that employed bilinear forms.

On an implementation level, the strong form can be implemented by taking $v = \phi_k$ for every $|k| \leq N$. For a fixed k , this reads,

$$\left\langle \sum_{|\ell| \leq N} -(i\ell)^2 \hat{u}_\ell \phi_\ell, \phi_k \right\rangle + \langle u_N, \phi_k \rangle = \langle f, \phi_k \rangle.$$

This is *exactly*,

$$\widehat{D}_2 \hat{u} = \hat{f},$$

that is equivalent to the previous Fourier-Galerkin scheme we derived.

Caution: It is not always the case that the strong variational form and weak variational form coincide.

Beyond Fourier-Galerkin

It's worthwhile to mention how this error analysis might be generalized to another problem or setup:

- The crucial components were Céa's Lemma and fundamental approximation estimates
- There is little theoretical difference if the spatial domain is multidimensional, or if the basis functions are different.
- For variable coefficient problems, e.g.,

$$-\frac{d}{dx} \left((2 + \sin x) \frac{d}{dx} u(x) \right) + u(x) = f(x),$$

the main difficulty is computational: the term $(2 + \sin x)$ results in the need to compute more complicated inner products, e.g.,

$$\langle (2 + \sin x) \phi'_k(x), \phi'_\ell(x) \rangle.$$

This is a computational issue, but does not adversely affect theory.

- For non-periodic problems: the theoretical strategies are quite similar, and much of the marginal extra work involves fundamental approximation estimates for non-periodic basis functions
- Time-dependent problems require somewhat different theoretical strategies, but *stability* is a big piece of the puzzle.

Fourier-Collocation

The Fourier collocation strategy uses a very similar approach to the Galerkin case.

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's}$$

As with the Galerkin method, our trial space U is the *Fourier* space:

$$u \in \text{span}\{\phi_k(x)\}_{k \in \mathbb{Z}}, \quad \phi_k(x) = \frac{1}{\sqrt{2\pi}} e^{ikx}.$$

However, our test space V is formally defined by Dirac functions centered at the grid points:

$$V = \text{span}\{\delta_{x_m}\}_{m \in [2N+1]}, \quad x_m = \frac{2\pi(m-1)}{2N+1}.$$

Fourier-Collocation

The Fourier collocation strategy uses a very similar approach to the Galerkin case.

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's}$$

As with the Galerkin method, our trial space U is the *Fourier* space:

$$u \in \text{span}\{\phi_k(x)\}_{k \in \mathbb{Z}}, \quad \phi_k(x) = \frac{1}{\sqrt{2\pi}} e^{ikx}.$$

However, our test space V is formally defined by Dirac functions centered at the grid points:

$$V = \text{span}\{\delta_{x_m}\}_{m \in [2N+1]}, \quad x_m = \frac{2\pi(m-1)}{2N+1}.$$

I.e., our weighted residual method in the collocation setup is quite simple:

$$-u_N''(x_m) + u_N(x_m) = f(x_m), \quad m \in [2N+1].$$

Equivalently, we enforce the condition,

$$I_N \mathcal{R}(u_N) = 0,$$

which is an exact equality since zero residual on interpolation nodes implies that the interpolant is 0.

This method is quite informal as we simply enforce pointwise zero residual.

Implementation

Recall some Fourier notation:

$$u_N(x) = \sum_{|k| \leq N} \hat{u}_k \phi_k(x), \quad \hat{\mathbf{u}} = \begin{pmatrix} \hat{u}_{-N} \\ \vdots \\ \hat{u}_N \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} u(x_1) \\ \vdots \\ u(x_M) \end{pmatrix}.$$

$$\tilde{\mathbf{u}} = \tilde{\mathbf{V}}^* \mathbf{u}, \quad \mathbf{u} = \frac{M}{2\pi} \tilde{\mathbf{V}} \tilde{\mathbf{u}}.$$

There are two computational options for implementing collocation:

- Store/solve for $\hat{\mathbf{u}}$ as degrees of freedom
- Store/solve for \mathbf{u} as degrees of freedom.

The latter, being more consistent with the idea of collocation, is more standard.

Implementation

Recall some Fourier notation:

$$u_N(x) = \sum_{|k| \leq N} \hat{u}_k \phi_k(x), \quad \hat{\mathbf{u}} = \begin{pmatrix} \hat{u}_{-N} \\ \vdots \\ \hat{u}_N \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} u(x_1) \\ \vdots \\ u(x_M) \end{pmatrix}.$$

$$\tilde{\mathbf{u}} = \tilde{\mathbf{V}}^* \mathbf{u}, \quad \mathbf{u} = \frac{M}{2\pi} \tilde{\mathbf{V}} \tilde{\mathbf{u}}.$$

There are two computational options for implementing collocation:

- Store/solve for $\hat{\mathbf{u}}$ as degrees of freedom
- Store/solve for \mathbf{u} as degrees of freedom.

The latter, being more consistent with the idea of collocation, is more standard.

To implement such a scheme, the difficult part is really to compute $u_N''(x_m)$. I.e., we must differentiate twice and evaluate at x_m . This can be accomplished by:

1. Transform \mathbf{u} to $\hat{\mathbf{u}}$. (Apply $\tilde{\mathbf{V}}^*$.)
2. Twice differentiate $u_N = \sum \hat{u}_k \phi_k$. (Apply $\text{diag}((-N)^2, \dots, N^2)$.)
3. Transform $\hat{\mathbf{u}}$ back to \mathbf{u} . (Apply $\frac{M}{2\pi} \tilde{\mathbf{V}}$.)

Differentiation matrices

This allows us to create a *nodal* equivalent of the continuous operation,

$$u_N(x) \mapsto u_N''(x)$$

This operation, with nodal degrees of freedom, is equivalent to applying the matrix,

$$\tilde{\mathbf{D}}_2 := \frac{M}{2\pi} \tilde{\mathbf{V}} \begin{pmatrix} (-N)^2 & & & \\ & (-N+1)^2 & & \\ & & \ddots & \\ & & & (N)^2 \end{pmatrix} \tilde{\mathbf{V}}^*.$$

Differentiation matrices

This allows us to create a *nodal* equivalent of the continuous operation,

$$u_N(x) \mapsto u_N''(x)$$

This operation, with nodal degrees of freedom, is equivalent to applying the matrix,

$$\tilde{\mathbf{D}}_2 := \frac{M}{2\pi} \tilde{\mathbf{V}} \begin{pmatrix} (-N)^2 & & & & \\ & (-N+1)^2 & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & (N)^2 \end{pmatrix} \tilde{\mathbf{V}}^*.$$

This matrix is dense and is equivalent to a finite difference matrix whose stencil spans the entire grid. (Except we approximate with Fourier Series and not polynomials.)

With some abuse of notation, let $I_N \mathbf{u}$ denote the interpolant (an element of V_N) corresponding to the unique element whose grid values are \mathbf{u} . Then,

$$I_N \tilde{\mathbf{D}}_2 \mathbf{u} = (I_N \mathbf{u})'',$$

i.e., $\tilde{\mathbf{D}}_2$ accomplishes *exact* second differentiation for elements from V_N .

Fourier collocation scheme

The Fourier collocation scheme associated to

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's,}$$

then reads,

$$\left(-\tilde{\mathbf{D}}_2 + \mathbf{I}\right) \mathbf{u} = \mathbf{f}, \quad \mathbf{f} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_M) \end{pmatrix},$$

which is a linear system that can be solved for \mathbf{u} , which uniquely identifies $I_N \mathbf{u} \in V_N$.

Note here that $\left(-\tilde{\mathbf{D}}_2 + \mathbf{I}\right)$ is a dense matrix, so some advantages from the Fourier-Galerkin setting are missing.

Fourier collocation scheme

The Fourier collocation scheme associated to

$$-u''(x) + u(x) = f(x), \quad \text{periodic BC's,}$$

then reads,

$$\left(-\tilde{\mathbf{D}}_2 + \mathbf{I}\right) \mathbf{u} = \mathbf{f}, \quad \mathbf{f} = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_M) \end{pmatrix},$$

which is a linear system that can be solved for \mathbf{u} , which uniquely identifies $I_N \mathbf{u} \in V_N$.

Note here that $\left(-\tilde{\mathbf{D}}_2 + \mathbf{I}\right)$ is a dense matrix, so some advantages from the Fourier-Galerkin setting are missing.

What about error estimates? For this particularly simple equation and for *smooth* solutions u_N^1 , note that the Fourier-Galerkin method for $u_{N,G}$ and the Fourier-Collocation method for $u_{N,C}$ are given by, respectively, is given by:

$$P_N \left(-u''_{N,G} + u_{N,G}\right) = P_N f, \quad I_N \left(-u''_{N,C} + u_{N,C}\right) = I_N f,$$

where P_N is the L^2 -orthogonal projector onto V_N .

¹I.e., u_N smooth enough so that $\langle -u''_N, v \rangle = \langle u'_N, v' \rangle$, which is always true if $u_N, v \in V_N$.

Collocation error estimates, I

$$P_N \left(-u''_{N,G} + u_{N,G} \right) = P_N f, \quad I_N \left(-u''_{N,C} + u_{N,C} \right) = I_N f,$$

Since $-u''_N + u \in V_N$, then

$$I_N \left(-u''_{N,C} + u_{N,C} \right) = P_N \left(-u''_{N,C} + u_{N,C} \right),$$

and therefore the difference $\Delta u_N := u_{N,G} - u_{N,C}$ satisfies,

$$P_N \left(-\Delta u''_N + \Delta u_N \right) = A_N f, \quad = P_N A_N f$$

where $A_N f = P_N f - I_N f$ is the aliasing error for f .

Collocation error estimates, I

$$P_N \left(-u''_{N,G} + u_{N,G} \right) = P_N f, \quad I_N \left(-u''_{N,C} + u_{N,C} \right) = I_N f,$$

Since $-u''_N + u \in V_N$, then

$$I_N \left(-u''_{N,C} + u_{N,C} \right) = P_N \left(-u''_{N,C} + u_{N,C} \right),$$

and therefore the difference $\Delta u_N := u_{N,G} - u_{N,C}$ satisfies,

$$P_N \left(-\Delta u''_N + \Delta u_N \right) = A_N f,$$

where $A_N f = P_N f - I_N f$ is the aliasing error for f .

Hence, the error between $u_{N,C}$ and $u_{N,G}$ is given by the solution to $-u'' + u = g$, where g is the aliasing error.

At the discrete level, this is,

$$\left(\widehat{D}_2 + \mathbf{I} \right) \widehat{\Delta u_N} = \left(\widehat{f} - \widetilde{f} \right),$$

where \widetilde{f} are the expansion coefficients for f computed using interpolation I_N .

Collocation error estimates, II

Therefore, the norm of Δu_N satisfies,

$$\begin{aligned}\|\Delta u_N\|_2 &\leq \left\| \left(\widehat{\mathbf{D}}_2 + \mathbf{I} \right)^{-1} \right\|_2 \|\widehat{\mathbf{f}} - \widetilde{\mathbf{f}}\|_2 \\ &= \left\| \left(\widehat{\mathbf{D}}_2 + \mathbf{I} \right)^{-1} \right\|_2 \|A_N f\|_{L^2} \\ &\leq \|A_N f\|_{L^2} \lesssim N^{-r}.\end{aligned}$$

where we have assumed $f \in H_p^r$.

Collocation error estimates, II

Therefore, the norm of Δu_N satisfies,

$$\begin{aligned}\|\Delta u_N\|_2 &\leq \left\| \left(\widehat{\mathbf{D}}_2 + \mathbf{I} \right)^{-1} \right\|_2 \|\widehat{\mathbf{f}} - \widetilde{\mathbf{f}}\|_2 \\ &= \left\| \left(\widehat{\mathbf{D}}_2 + \mathbf{I} \right)^{-1} \right\|_2 \|A_N f\|_{L^2} \\ &\leq \|A_N f\|_{L^2} \lesssim N^{-r}.\end{aligned}$$

where we have assumed $f \in H_p^r$.

Chaining this with our H_p^1 -norm estimate for $u - P_N u$ yields,

$$f \in H_p^r \implies \|u - u_{N,C}\|_{L^2} \lesssim N^{-r},$$

where we have used the fact that $f \in H_p^r$ implies that $u \in H_p^{r+2}$.

A finer estimate can be constructed with sharper treatment of the last inequality in the estimate for $\|\Delta u_N\|_2$.

(I.e., the estimate we've derived is not optimal, but one can derive an optimal estimate.)

One can also construct an estimate in $V = H_p^1$.

Galerkin vs collocation, I

For the particular example we've considered:

In practice, both the Galerkin and collocation schemes **behave quite similarly for smooth u** .

Both schemes can utilize the FFT: For the Galerkin method this is straightforward. For collocation, the application $v \mapsto \widehat{D}_2$ can be accomplished by sandwiching application of the (diagonal) matrix \widehat{D}_2 with the forward/inverse FFT.

If quadrature/interpolation is used to approximate \widehat{f} in the Galerkin scheme, then **the schemes produce identical solutions**.

Galerkin vs collocation, I

For the particular example we've considered:

In practice, both the Galerkin and collocation schemes **behave quite similarly for smooth u** .

Both schemes can utilize the FFT: For the Galerkin method this is straightforward. For collocation, the application $v \mapsto \widehat{D}_2$ can be accomplished by sandwiching application of the (diagonal) matrix \widehat{D}_2 with the forward/inverse FFT.

If quadrature/interpolation is used to approximate \widehat{f} in the Galerkin scheme, then **the schemes produce identical solutions**.

The property is due to the equation we have considered. For example, if we instead consider,

$$-u''(x) + (2 + \sin x)u(x) = f(x),$$

then the Galerkin/collocation discretizations for the left hand side are *not* the same, and hence one must augment the previous analysis to estimate the discrepancy.

In general, it's easier to prove estimates for Galerkin methods.

Galerkin vs collocation, II

Collocation is more conceptually easy to implement, although constructing the differentiation matrix is somewhat nontrivial, even if one uses the FFT.

Nevertheless, an example such as,

$$-u''(x) + \frac{u(x)}{2 + \sin x} = f(x),$$

$$P_N \left[\frac{u_N(x)}{2 + \sin x} \right]$$

shows when collocation is more straightforward than Galerkin.

$$P_N \left[\frac{\phi_k(x)}{2 + \sin x} \right] = ?$$

Consider, in particular, when there are nonlinear terms in the equation: Galerkin methods can be difficult to implement.

In constrast, the collocation discretization of the above corresponds to

$$-u_N''(x_m) + \frac{u_N(x_m)}{2 + \sin x_m} = f(x_m), \quad m \in [M],$$

which is straightforward to implement as it is linear in u_N .

Galerkin vs collocation, III

One can borrow the interpolation idea from collocation to aid in approximating Galerkin quantities. I.e., the Galerkin scheme requires us to enforce,

$$P_N \left(-u_N''(x) + \frac{u_N(x)}{2 + \sin x} \right) = P_N f(x),$$

and computing $P_N(u_N/(2 + \sin x))$ is the difficult part. However, we could approximate this by:

$$P_N (-u_N''(x)) + I_N \left(\frac{u_N(x)}{2 + \sin x} \right) = P_N f(x),$$

which is easy to implement, as $I_N u_N/(2 + \sin x)$ is computed by evaluating at the gridpoints x_m and then taking a DFT.

This type of approximation of projections (in particular of nonlinear terms) is called a **pseudospectral** approximation.

References I

-  Canuto, Claudio et al. (2011). *Spectral Methods: Fundamentals in Single Domains*. 1st ed. 2006. Corr. 4th printing 2010 edition. Berlin ; New York: Springer. ISBN: 978-3-540-30725-9.
-  Shen, Jie, Tao Tang, and Li-Lian Wang (2011). *Spectral Methods: Algorithms, Analysis and Applications*. Springer Science & Business Media. ISBN: 978-3-540-71041-7.