

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF UTAH
Analysis of Numerical Methods, II
MATH 6620 – Section 001 – Spring 2024
Homework 5 Solutions
Fourier Series and time-dependent problems

Due Friday, April 5, 2024

Submit your solutions online through Gradescope.

1. (Hyperbolic systems)
 - a. A linear PDE,

$$\mathbf{u}_t + \mathbf{A}\mathbf{u}_x = 0,$$

governing the vector-valued solution $\mathbf{u}(x, t) \in \mathbb{R}^m$, is called *hyperbolic* if \mathbf{A} is diagonalizable and $\lambda(\mathbf{A}) \subset \mathbb{R}$, with $\lambda(\mathbf{A})$ denoting the spectrum of \mathbf{A} . Assuming this PDE is hyperbolic, and using an equispaced grid in time and space, derive an implementable version of the *upwind scheme* for this problem, where \mathbf{u}_t is discretized as $D^+\mathbf{u}_j^n$. You may ignore boundary conditions.

- b. Derive and state a CFL condition for your upwind scheme.
- c. Consider a *nonlinear* PDE,

$$\mathbf{u}_t + \mathbf{f}(\mathbf{u})_x = 0,$$

where $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$, and again we ignore boundary conditions. Propose a definition of a *hyperbolic PDE* in this context, and identify an upwind-like numerical scheme on an equidistant time and space mesh. (Again, discretize $\mathbf{u}_t \approx D^+\mathbf{u}_j^n$.) What is the CFL condition for your scheme?

Solution:

- a. By assumption, we have,

$$\mathbf{A} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1},$$

where $\mathbf{\Lambda}$ is a diagonal matrix with real values $\lambda_1, \dots, \lambda_m$ on the diagonal. We reformulate the PDE through the transformation,

$$\mathbf{w} := \mathbf{V}^{-1}\mathbf{u},$$

and so by multiplying the PDE by \mathbf{V}^{-1} , we obtain,

$$\mathbf{w}_t + \mathbf{\Lambda}\mathbf{w}_x = \mathbf{0},$$

which is an uncoupled system of PDEs of the form,

$$(w_i)_t + \lambda_i (w_i)_x = 0, \quad \mathbf{w}(x, t) = (w_1(x, t), \dots, w_m(x, t))^T.$$

An upwind scheme for each of these scalar PDEs would be,

$$D^+w_{i,j}^n + \lambda_i D_{\sigma_i} w_{i,j}^n = 0, \quad w_{i,j}^n \approx w_i(x_j, t_n),$$

where x_j is an equidistant grid in space with mesh spacing $h > 0$, and t_n is an equidistant grid in time with mesh spacing $k > 0$. The sign σ_i is determined by the sign of λ_i in the typical upwind manner for scalar, constant wavespeed problems:

$$\sigma_i := -\text{sign}(\lambda_j), \implies D_{\sigma_i} = \begin{cases} D_-, & \lambda_i \geq 0 \text{ (rightward moving wave)} \\ D_+, & \lambda_i < 0 \text{ (leftward moving wave)} \end{cases}$$

where it doesn't really matter how σ_i is defined when $\lambda_i = 0$. Then define the size- m vector of differences as,

$$D_{\boldsymbol{\sigma}} \mathbf{w}_j^n := \begin{pmatrix} D_{\sigma_1} w_{1,j}^n \\ D_{\sigma_2} w_{2,j}^n \\ \vdots \\ D_{\sigma_m} w_{1,m}^n \end{pmatrix}.$$

With this definition, then upwinding is given by,

$$D^+ \mathbf{w}_j^n + \Lambda D_{\boldsymbol{\sigma}} \mathbf{w}_j^n = 0$$

We translate this to \mathbf{u} space by multiplying by \mathbf{V} and replacing \mathbf{w}_j^n with its \mathbf{u} -representation:

$$D^+ \mathbf{u}_j^n + \mathbf{V} (\Lambda D_{\boldsymbol{\sigma}} \mathbf{V}^{-1} \mathbf{u}_j^n) = 0,$$

or more explicitly,

$$\mathbf{u}_j^{n+1} = \mathbf{u}_j^n - k \mathbf{V} (\Lambda D_{\boldsymbol{\sigma}} \mathbf{V}^{-1} \mathbf{u}_j^n)$$

- b. The scheme from the previous part is essentially just an upwind scheme in \mathbf{w} -space, so the CFL condition is the same as the CFL condition for \mathbf{w} space. For every $i \in [m]$, the CFL condition for the scheme evolving $w_{i,j}^n$ is,

$$k \leq \frac{h}{|\lambda_i|}$$

Hence, the CFL restriction for the whole system is the strictest condition for every i :

$$k \leq \min_{i \in [m]} \frac{h}{|\lambda_i|} = \frac{h}{\rho(\mathbf{A})},$$

where $\rho(\mathbf{A})$ is the spectral radius of \mathbf{A} :

$$\rho(\mathbf{A}) := \max_{i \in [m]} |\lambda_i|.$$

- c. In the linear case, we have,

$$\mathbf{f}(\mathbf{u}) = \mathbf{A}\mathbf{u},$$

and the definition for hyperbolicity was that \mathbf{A} is diagonalizable with spectrum $\lambda(\mathbf{A})$ satisfying,

$$\lambda(\mathbf{A}) \subset \mathbb{R}.$$

Restating this in a way that can generalize to the nonlinear case, the hyperbolicity in the linear case meant that $\mathbf{A} = \frac{\partial \mathbf{f}}{\partial \mathbf{u}}$ was diagonalizable and satisfies,

$$\lambda \left(\frac{\partial \mathbf{f}}{\partial \mathbf{u}} \right) = \lambda(\mathbf{A}) \subset \mathbb{R}.$$

This now serves as a potential definition: we might say that our nonlinear PDE is hyperbolic at a state \mathbf{u} if the Jacobian matrix $\frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{u})$ is a diagonalizable matrix, with,

$$\lambda \left(\frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{u}) \right) \subset \mathbb{R}.$$

Another way of motivating this definition is to rewrite the PDE by using the chain rule as $\mathbf{u}_t + \frac{\partial \mathbf{f}}{\partial \mathbf{u}} \mathbf{u}_x = \mathbf{0}$, so that the Jacobian matrix is the analog of the wavespeed matrix \mathbf{A} , and hence the Jacobian should be diagonalizable with real eigenvalues for every relevant \mathbf{u} . (E.g., every \mathbf{u} along a space-time solution trajectory.) Note that this condition seems rather strong, as we require such “real” diagonalizability for a relatively large class of vectors \mathbf{u} . Nevertheless, this gives us one way to identify a scheme: if the nonlinear system satisfies our definition of hyperbolicity, then we have for every space index j and every time index n ,

$$\frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{u}_j^n) = \mathbf{V}_j^n \mathbf{\Lambda}_j^n (\mathbf{V}_j^n)^{-1},$$

where we use sub- and superscripts to emphasize that quantities depend on the space and time index, respectively. (And in particular we must diagonalize \mathbf{f} at every time and space point). With this, we may transform the system and perform a scheme that is at least “locally” upwinding:

$$\mathbf{w}_j^n := (\mathbf{V}_j^n)^{-1} \implies D^+ \mathbf{w}_j^n + \mathbf{\Lambda}_j^n D \sigma_j^n \mathbf{w}_j^n = 0,$$

where $\mathbf{\Lambda}_j^n$ has diagonal entries $(\lambda_j^n)_1, \dots, (\lambda_j^n)_m$, and σ_j^n is given by,

$$(\sigma_j^n)_i = -\text{sign}(\lambda_j^n)_i.$$

Then transforming back into \mathbf{u} space and expanding the D^+ operator, we have,

$$\mathbf{u}_j^{n+1} = \mathbf{u}_j^n - k \mathbf{V}_j^n \left(\mathbf{\Lambda}_j^n D \sigma_j^n (\mathbf{V}_j^n)^{-1} \mathbf{u}_j^n \right)$$

Under the assumption that characteristics travel at constant speed determined by the Jacobian at \mathbf{u}_j^n , then the CFL condition is given by the assumption that the (j, n) -local characteristics lie within the numerical domain of dependence, i.e.,

$$k \leq \frac{h}{\max_{i \in [m]} (\lambda_j^n)_i} = \frac{h}{\rho \left(\frac{\partial \mathbf{f}}{\partial \mathbf{u}}(\mathbf{u}_j^n) \right)}$$

2. (Hyperbolic systems in two spatial dimensions) Consider a linear system of PDEs of the form,

$$\mathbf{u}_t + \mathbf{A} \mathbf{u}_x + \mathbf{B} \mathbf{u}_y = 0,$$

for the unknown $\mathbf{u}(x, y, t) \in \mathbb{R}^m$, where $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{m \times m}$ are given. Assume this system is *hyperbolic*, meaning that for every $\alpha, \beta \in \mathbb{R}$, then $\alpha \mathbf{A} + \beta \mathbf{B}$ is diagonalizable with $\lambda(\alpha \mathbf{A} + \beta \mathbf{B}) \subset \mathbb{R}$. We will use the abbreviation $\mathbf{x} = (x, y)^T \in \mathbb{R}^2$.

- a. Determine plane wave solutions: let $\mathbf{n} \in \mathbb{R}^2$ satisfying $\|\mathbf{n}\|_2 = 1$ be a given vector, and assume some initial data $\mathbf{u}(x, y, 0) = \mathbf{u}_0(\mathbf{n} \cdot \mathbf{x})$. (Note that \mathbf{u}_0 is a function of a scalar input.) Ignoring boundary conditions, determine the solution $\mathbf{u}(x, y, t)$.
- b. Assume further that \mathbf{A} and \mathbf{B} are symmetric matrices. Consider discretizing this PDE on an equidistant spatial grid on a square: $h_x > 0$ and $h_y > 0$ are the grid spacing in the x - and y -directions, respectively. Again, assume the discretization $\mathbf{u}_t \approx D^+ \mathbf{u}_{i,j}^n$, where $\mathbf{u}_{i,j}^n \approx \mathbf{u}(x_i, y_j, t_n)$. Using plane waves as motivation, derive a CFL condition (a condition on the timestep k) for such a discretization. You may assume that spatial derivatives are approximated using the stencil involving $u_{i\pm 1, j\pm 1}^n$ and that the numerical domain of dependence is the convex hull of the stencil points (as in the one-dimensional case).

Solution:

- a. Given $\mathbf{n} = (n_x, n_y)^T$ with $n_x^2 + n_y^2 = 1$, define $r := \mathbf{n} \cdot \mathbf{x} = xn_x + yn_y$. Thus, the initial data is $\mathbf{u}(x, y, 0) = \mathbf{u}_0(r)$. This suggests we should look for solutions of the form $\mathbf{u}(x, y, t) =: \mathbf{w}(xn_x + yn_y, t) = \mathbf{w}(r, t)$. In order to determine a PDE for \mathbf{w} , note that,

$$\begin{aligned} \frac{\partial}{\partial t} \mathbf{u}(x, y, t) &= \frac{\partial}{\partial t} \mathbf{w}(r, t), \\ \frac{\partial}{\partial x} \mathbf{u}(x, y, t) &= \frac{\partial r}{\partial x} \frac{\partial}{\partial r} \mathbf{u}(x, y, t) = n_x \frac{\partial}{\partial r} \mathbf{u}(x, y, t) \\ \frac{\partial}{\partial y} \mathbf{u}(x, y, t) &= \frac{\partial r}{\partial y} \frac{\partial}{\partial r} \mathbf{u}(x, y, t) = n_y \frac{\partial}{\partial r} \mathbf{u}(x, y, t) \end{aligned}$$

Therefore,

$$\mathbf{u}_t + \mathbf{A}\mathbf{u}_x + \mathbf{B}\mathbf{u}_y = 0 \quad \longrightarrow \quad \mathbf{w}_t + (n_x \mathbf{A} + n_y \mathbf{B}) \mathbf{w}_r = 0.$$

Since by assumption $n_x \mathbf{A} + n_y \mathbf{B}$ is diagonalizable with real spectrum,

$$n_x \mathbf{A} + n_y \mathbf{B} = \mathbf{V}_n \mathbf{\Lambda}_n \mathbf{V}_n^{-1}, \quad \mathbf{\Lambda}_n = \text{diag}(\lambda_{n,1}, \dots, \lambda_{n,m}), \quad \lambda_{n,q} \in \mathbb{R} \quad \forall q \in [m],$$

where we have used \mathbf{n} subscripts to emphasize quantities that depend on \mathbf{n} . Then the PDE for \mathbf{w} is a system of linear hyperbolic PDEs in a scalar space dimension r . The solution is given by linearly transforming \mathbf{w} into the variable $\mathbf{W} := \mathbf{V}_n^{-1} \mathbf{w}$ to decouple the system:

$$\mathbf{W}_t + \mathbf{\Lambda}_n \mathbf{W}_r = 0, \quad \mathbf{W}(r, 0) = \mathbf{V}_n^{-1} \mathbf{w}(r, 0) = \mathbf{V}_n^{-1} \mathbf{u}_0(r).$$

Hence, the exact solution for \mathbf{W} and hence \mathbf{u} is given by,

$$\begin{aligned} \mathbf{W}(r, t) &= \mathbf{W}(r - \lambda_n t, 0) := \begin{pmatrix} W_1(r - \lambda_{n,1} t, 0) \\ \vdots \\ W_m(r - \lambda_{n,m} t, 0) \end{pmatrix} \\ \mathbf{u}(x, y, t) &= \mathbf{V}_n \mathbf{W}(xn_x + yn_y - \lambda_n t, 0). \end{aligned}$$

Through this argument, we have constructed a particular solution \mathbf{u} to the PDE (that satisfies the initial data). Note that more general solutions, e.g., $u = u(r, s, t)$ with nontrivial s dependence, with $s = -n_y x + n_x y$, cannot exist due to well-posedness and linearity of the PDE. If such alternative solutions did exist, then the difference between the s -dependent and s -independent solution above would satisfy the same PDE but with zero initial data, meaning that the solution difference would be the zero function.

- b. Note that for fixed \mathbf{n} , the solution to the above system for \mathbf{u} (or \mathbf{W}) involves waves traveling in \mathbb{R}^2 in the direction \mathbf{n} at speeds,

$$\lambda_{\mathbf{n},1}, \dots, \lambda_{\mathbf{n},m}.$$

For the one-directional dimension \mathbf{n} , we consider the origin-centered numerical domain of dependence to be the convex hull of neighboring points:

$$\text{numerical domain of dependence} = \text{conv} (h_x(1, 0)^T, h_x(-1, 0)^T, h_y(0, 1)^T, h_y(0, -1)^T),$$

which is an origin-centered rotated quadrilateral. The distance from the origin to the boundary of this set along the direction \mathbf{n} is given by,

$$h_r := \frac{1}{\frac{|n_x|}{h_x} + \frac{|n_y|}{h_y}} = \frac{h_x h_y}{h_y |n_x| + h_x |n_y|}$$

Hence, in the one-dimensional direction \mathbf{n} , this problem has the following CFL restriction:

$$k \leq \min_{q \in [m]} \frac{h_r}{|\lambda_{\mathbf{n},q}|} = \frac{h_r}{\rho(n_x \mathbf{A} + n_y \mathbf{B})}.$$

While this is the CFL restriction considering waves along a direction \mathbf{n} , the more practical restriction is an \mathbf{n} -independent condition. I.e., we should seek a condition without *a priori* knowledge of \mathbf{n} . Then formally this would be the condition,

$$k \leq \sup_{\|\mathbf{n}\|_2=1} \frac{h_r}{\rho(n_x \mathbf{A} + n_y \mathbf{B})} \quad (1)$$

In order to derive a computable condition from this, we exploit the further assumptions that \mathbf{A} and \mathbf{B} are symmetric. We claim that under this extra assumption,

$$\rho(n_x \mathbf{A} + n_y \mathbf{B}) \leq |n_x| \rho(\mathbf{A}) + |n_y| \rho(\mathbf{B}). \quad (2)$$

To prove this, note that since \mathbf{A} and \mathbf{B} are symmetric, then their spectra are characterized by the Rayleigh quotient. E.g., for \mathbf{A} :

$$\lambda_{\min}(\mathbf{A}) = \min_{\mathbf{x} \neq \mathbf{0}} R_{\mathbf{A}}(\mathbf{x}), \quad \lambda_{\max}(\mathbf{A}) = \max_{\mathbf{x} \neq \mathbf{0}} R_{\mathbf{A}}(\mathbf{x}), \quad \rho(\mathbf{A}) = \max_{\mathbf{x} \neq \mathbf{0}} |R_{\mathbf{A}}(\mathbf{x})|,$$

with,

$$R_{\mathbf{A}}(\mathbf{x}) := \frac{\mathbf{x}^T \mathbf{A} \mathbf{x}}{\mathbf{x}^T \mathbf{x}}.$$

Note in particular that $R_{\mathbf{A}}(\cdot)$ is linear in \mathbf{A} . Using this characterization, then,

$$\begin{aligned} \rho(n_x \mathbf{A} + n_y \mathbf{B}) &= \max_{\mathbf{x} \neq \mathbf{0}} |R_{n_x \mathbf{A} + n_y \mathbf{B}}(\mathbf{x})| \\ &= \max_{\mathbf{x} \neq \mathbf{0}} |n_x R_{\mathbf{A}}(\mathbf{x}) + n_y R_{\mathbf{B}}(\mathbf{x})| \\ &\leq |n_x| \max_{\mathbf{x} \neq \mathbf{0}} |R_{\mathbf{A}}(\mathbf{x})| + |n_y| \max_{\mathbf{x} \neq \mathbf{0}} |R_{\mathbf{B}}(\mathbf{x})| \\ &= |n_x| \rho(\mathbf{A}) + |n_y| \rho(\mathbf{B}), \end{aligned}$$

as claimed in (2). To see how (2) assists us with coming up with a more transparent CFL condition, note that we have just shown that,

$$\frac{h_r}{|n_x|\rho(\mathbf{A}) + |n_y|\rho(\mathbf{B})} \leq \frac{h_r}{|\rho(n_x\mathbf{A} + n_y\mathbf{B})|},$$

and hence is we can assure that k is smaller than the left-hand side above for all \mathbf{n} , then it will satisfy our CFL condition (1). This is enough to create a servicable coarse bound, since,

$$\begin{aligned} \frac{1}{h_r} (|n_x|\rho(\mathbf{A}) + |n_y|\rho(\mathbf{B})) &= \left(\frac{|n_x|}{h_x} + \frac{|n_y|}{h_y} \right) (|n_x|\rho(\mathbf{A}) + |n_y|\rho(\mathbf{B})) \\ &\leq \left(\frac{1}{h_x} + \frac{1}{h_y} \right) (\rho(\mathbf{A}) + \rho(\mathbf{B})), \end{aligned}$$

so that a sufficient CFL condition reads,

$$k \leq \frac{h_x h_y}{(h_x + h_y) (\rho(\mathbf{A}) + \rho(\mathbf{B}))},$$

where this bound ensures (1). (One may of course create a tighter bound by optimizing over all \mathbf{n} , but this produces only a slightly better bound.) Note that this bound has all the right general behavior for $h_x = h_y$ and/or $\rho(\mathbf{A}) = \rho(\mathbf{B})$.

3. (Finite differences for non-smooth problems) Consider Burgers' equation:

$$u_t + f(u)_x = 0, \quad f(u) = \frac{u^2}{2}, \quad (x, t) \in [-\pi, \pi] \times (0, T].$$

where we will take $T = 1.0$. Supplement this PDE with the boundary conditions,

$$u(\pm\pi, t) = u(\pm\pi, 0),$$

where the initial condition function $u(\cdot, 0)$ will be specified below. Note that for smooth u , then $f(u)_x = uu_x$. Based on this observation, and using an equidistant grid in both space and time, we will consider two schemes for this PDE:

$$\text{Scheme A: } D^+ u_j^n + D_0 f(u_j^n) = 0,$$

$$\text{Scheme B: } D^+ u_j^n + u_j^n D_0 u_j^n = 0.$$

Numerically test these schemes for solving the PDE up to time $t = T$ with the following three initial data:

$$u(x, 0) = u_1(x) = -\sin(x). \tag{3a}$$

$$u(x, 0) = u_2(x) = \begin{cases} 1, & x \leq 0 \\ 0, & x > 0 \end{cases} \tag{3b}$$

$$u(x, 0) = u_3(x) = \begin{cases} 0, & x \leq 0 \\ 1, & x > 0 \end{cases} \tag{3c}$$

Based on the experiments above, which scheme would you prefer to use for each example, and in general? (Feel free to try other schemes as well, e.g., upwind versions of Schemes A and B as identified in problem 1c.)

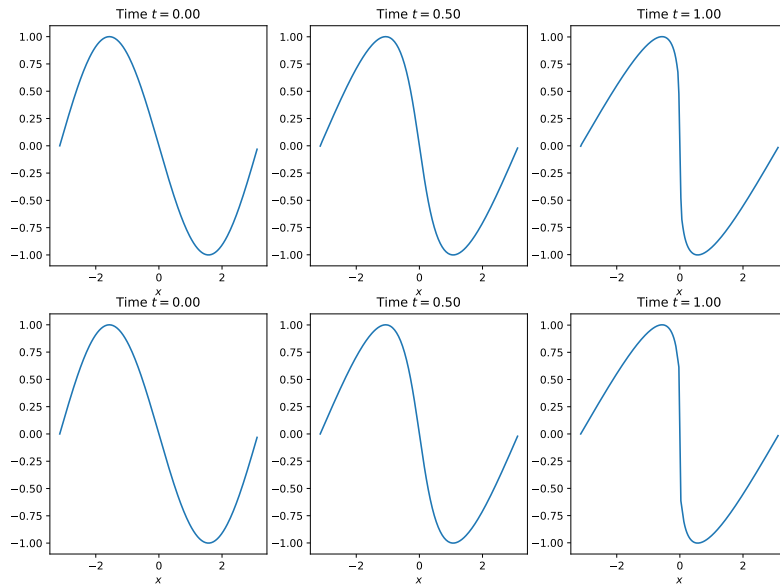


Figure 1: Problem 3: Scheme A (top) and Scheme B (bottom) for the initial condition (3a) to Burgers' equation. The left, center, and right columns correspond to times $t = 0, 0.5,$ and 1.0 .

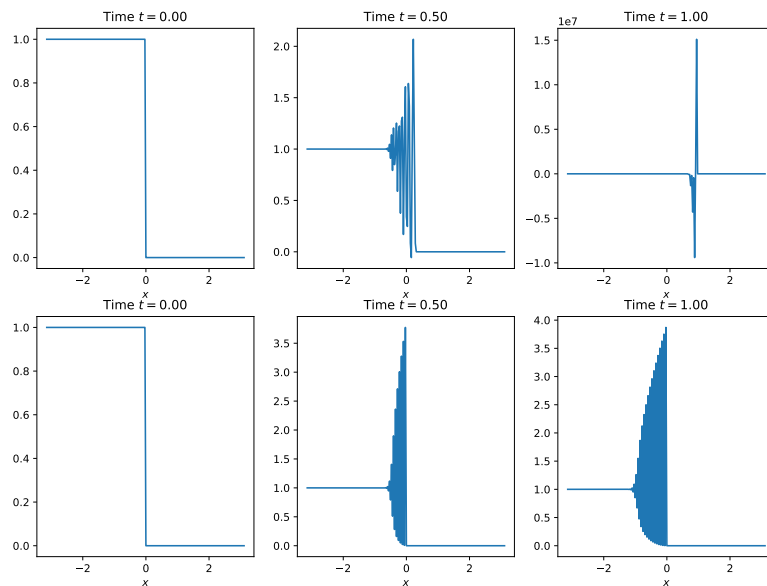


Figure 2: Problem 3: Scheme A (top) and Scheme B (bottom) for the initial condition (3b) to Burgers' equation. The left, center, and right columns correspond to times $t = 0, 0.5,$ and 1.0 .

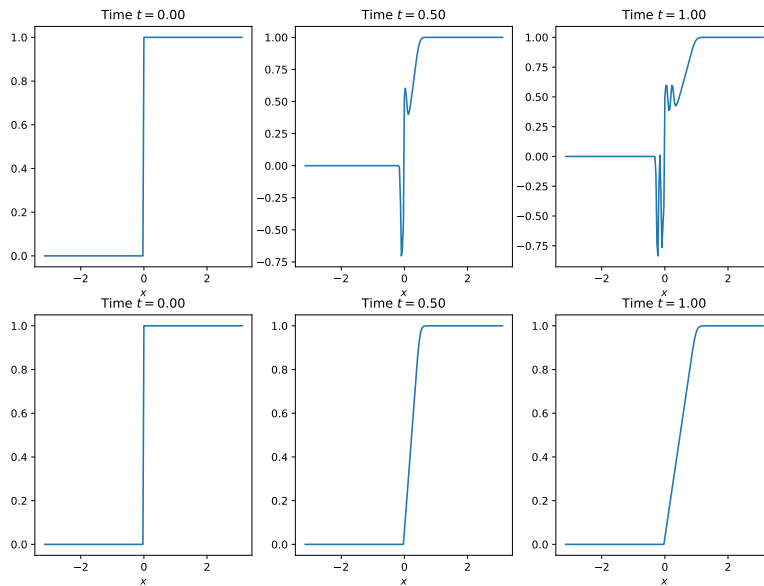


Figure 3: Problem 3: Scheme A (top) and Scheme B (bottom) for the initial condition (3c) to Burgers' equation. The left, center, and right columns correspond to times $t = 0, 0.5,$ and 1.0 .

Solution: Using the schemes described, we simulate the solution up to time $T = 1$ using $M = 200$ equispaced points in space, and $N = 200$ timesteps. (I.e., $h = 2\pi/M$ and $k = 1/N$.) The numerical results for initial condition (3a) are presented in figure 1, for (3b) in figure 2, and for (3c) in figure 3.

The first major observation is that for initial condition (3a), both schemes perform comparably. This is perhaps not surprising since both solutions appear to be relatively smooth functions and these two schemes are essentially identical so long as $f(u)$ is smooth. However, for both discontinuous initial conditions (3b) and (3c), it's not quite clear which scheme is better. Both schemes produce seemingly unphysical oscillations for initial condition (3b), and both produce a somewhat reasonable-looking solution for initial condition (3c). Thus, these schemes produce significantly different behavior for non-smooth initial data, and it's not clear which one is better.

In fact, neither of these schemes is bulletproof for these examples: For initial condition (3c), it turns out that Scheme B produces a numerical solution that is at least a reasonable approximation to the "true" solution. In contrast, Scheme A produces something that is stable, but "very" incorrect. However, when we look at initial condition (3b), both schemes produce rather bad oscillations, and Scheme A is (eventually) unstable. However, Scheme A for this case actually produces a solution that has some of the correct characteristics; in particular, the real solution to this problem is a moving discontinuity, and Scheme A produces an essentially correct movement of the discontinuity, even if it's accompanied by very incorrect oscillations and eventual instability. While Scheme B appears stable for this initial condition, the discontinuity does not move in the correct way, and for example at $t = 0.5$, the Scheme B solution is worse than the Scheme A solution. (Although again both are quite bad.) This highlights a major point: these two initial conditions are seemingly the same experiment, but in fact

the correctness of each scheme varies considerably between these cases. The major point here is that the methods we've developed are not quite suitable (at least without modification) to tackle nonlinear problems with non-smooth solutions.

4. (Fourier approximation)

- a. Prove that if $u \in H_p^s$ for some integer $s \geq 0$, then

$$\|u - P_N u\|_{L^2} \leq N^{-s} \|u\|_{H_p^s},$$

where the space H_p^s is defined on slide D13-S10 and P_N is defined in slide D13-S07.

- b. Confirm this behavior by numerically computing $\|u_j - P_N u_j\|_{L^2}$ as a function of N for each $j = 0, 1, 2, 3$. The functions u_j , $j \geq 0$, are defined as,

$$u_0(x) := \begin{cases} 1, & |x - \pi| < \frac{\pi}{2} \\ -1, & \text{else} \end{cases} \quad u_j(x) := c_j + \int_0^x u_{j-1}(y) dy \quad (j \geq 1),$$

where c_j is chosen so that u_j is a mean-0 function. Based on your numerical results what type of regularity (s) does u_j seem to have?

Solution:

- a. If $s = 0$, the desired statement is,

$$\|u - P_N u\|_{L^2} \leq \|u\|_{L^2},$$

which is true since $(I - P_N)$ is an orthogonal projector. Therefore, suppose $s \geq 1$. We have,

$$u \in H_p^s \implies u^{(j)} \in L^2, \quad 0 \leq j \leq s, \quad u^{(\ell)}(0) = u^{(\ell)}(2\pi), \quad 0 \leq \ell \leq s - 1.$$

Hence, for $j \leq s$, we define expansion coefficients for any $k \in \mathbb{Z}$,

$$\widehat{u}_k^{(j)} := \left\langle u^{(j)}(x), \phi_k(x) \right\rangle = \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} u^{(j)}(x) e^{-ikx} dx.$$

These coefficients satisfy Parseval's relation,

$$\|u^{(j)}\|_{L^2}^2 = \sum_{k \in \mathbb{Z}} |\widehat{u}_k^{(j)}|^2.$$

Through integration by parts, we find for any $1 \leq j \leq s$:

$$\begin{aligned} \widehat{u}_k^{(j)} &= \frac{1}{\sqrt{2\pi}} \int_0^{2\pi} u^{(j)}(x) e^{-ikx} dx \\ &\stackrel{\text{(IbP)}}{=} -\frac{ik}{\sqrt{2\pi}} e^{-ikx} u^{(j-1)}(x) \Big|_0^{2\pi} + \frac{ik}{\sqrt{2\pi}} \int_0^{2\pi} u^{(j-1)}(x) e^{-ikx} dx \\ &\stackrel{(*)}{=} \frac{ik}{\sqrt{2\pi}} \int_0^{2\pi} u^{(j-1)}(x) e^{-ikx} dx \\ &= ik \widehat{u}_k^{(j-1)}. \end{aligned}$$

where $(*)$ uses $u^{(j-1)}(0) = u^{(j-1)}(2\pi)$, which is true since $0 \leq j - 1 \leq s - 1$. Hence, by finite induction, for any $k \neq 0$:

$$\widehat{u}_k = \widehat{u}_k^{(0)} = \frac{\widehat{u}_k^{(s)}}{(ik)^s}. \tag{4}$$

Finally, for any $N \geq 1$:

$$\begin{aligned}
 \|u - P_N u\|_{L^2}^2 &= \left\| \sum_{|k| > N} \widehat{u}_k \phi_k \right\|_{L^2}^2 = \sum_{|k| > N} |\widehat{u}_k|^2 \\
 &\stackrel{(4)}{=} \sum_{|k| > N} \frac{1}{k^{2s}} |\widehat{u}_k^{(s)}|^2 \\
 &\leq N^{-2s} \sum_{|k| > N} |\widehat{u}_k^{(s)}|^2 \\
 &\leq N^{-2s} \sum_{k \in \mathbb{Z}} |\widehat{u}_k^{(s)}|^2 \\
 &= N^{-2s} \|u^{(s)}\|_{L^2}^2 \\
 &\leq N^{-2s} \sum_{j=0}^s \|u^{(j)}\|_{L^2}^2 = N^{-2s} \|u\|_{H_p^s}^2
 \end{aligned}$$

which proves the result.

- b. One can very well proceed to compute the requested norms by numerical computation (involving some approximation). Some options include computing projection coefficients approximately using dense quadrature, or even using simple $(2N + 1)$ -point interpolation since question 5 below establishes that interpolation will produce essentially the same asymptotic error. However, here is a way to compute the exact L^2 error norms: First we observe the following: if f is a given function, and $s \geq 1$, then

$$f \in L^2, \quad f(0) = f(2\pi), \quad \text{and } f' \in H_p^{s-1} \quad \implies f \in H_p^s,$$

which follows directly from the definition of H_p^s . Coming to the current problem, we see directly that $u_0 \in L^2 = H_p^0$. Furthermore, for $j \geq 1$:

$$u_j(2\pi) = c_j + \int_0^{2\pi} u_{j-1}(y) dy \stackrel{(*)}{=} c_j = u_j(0),$$

where $(*)$ uses that u_{j-1} is a mean-0 function. Furthermore, $u_j(x)$ by definition is a continuous function on the closed interval $[0, 2\pi]$. Therefore, it has a finite maximum inside this interval, and hence is an L^2 function. Hence, by finite induction, $u_j \in H_p^j$. Hence, the computations from part (a) allow us to conclude for $k \in \mathbb{Z} \setminus \{0\}$:

$$u_j \in H_p^j \quad \text{and} \quad u_j^{(j)}(x) = u_0(x) \quad \implies \quad \widehat{u}_{j,k} = \frac{\widehat{u}_{j,0}}{(ik)^j},$$

where $\widehat{u}_{j,k}$ is the ϕ_k -Fourier expansion coefficient for u_j . Therefore, we need only compute the Fourier coefficients for u_0 in order to determine the exact coefficients for all the remaining functions. (Since u_j is mean-0 for every j , then $\widehat{u}_{j,0} = 0$.) We compute the

coefficients for u_0 , $k \neq 0$, explicitly:

$$\begin{aligned} \langle u_0, \phi_k \rangle &= \frac{1}{\sqrt{2\pi}} \left[\int_{\pi/2}^{3\pi/2} e^{-ikx} dx - \int_{-\pi/2}^{\pi/2} e^{-ikx} dx \right] \\ &= -\frac{1}{ik\sqrt{2\pi}} \left[e^{-ikx} \Big|_{\pi/2}^{3\pi/2} - e^{-ikx} \Big|_{-\pi/2}^{\pi/2} \right] \\ &= -\frac{4}{\sqrt{2\pi}} \frac{\sin(k\pi/2)}{k} \\ &= \begin{cases} \frac{-4(-1)^{(k-1)/2}}{k\sqrt{2\pi}}, & k \text{ odd} \\ 0, & k \text{ even} \end{cases} \end{aligned}$$

Then this gives us the result:

$$\hat{u}_{j,k} = \begin{cases} \frac{-4(-1)^{(k-1)/2}}{i^j k^{j+1} \sqrt{2\pi}}, & k \text{ odd} \\ 0, & k \text{ even} \end{cases}$$

which holds even for $k = 0$. In order to compute the error norms, we will use the fact that,

$$\|u - P_N u\|_{L^2}^2 = \|u\|_{L^2}^2 - \sum_{|k| \leq N} |\hat{u}_k|^2,$$

which is the Pythagorean theorem $\|u\|^2 = \|P_N u\|^2 + \|(I - P_N)u\|^2$. Note that this furnishes an exact formula for the L^2 error if we know the L^2 norms of u_j along with a finite number of expansions coefficients (which we just computed above). In order to compute the exact L^2 norms, instead of taking antiderivatives to compute formulas for u_j , we utilize Parseval's identity:

$$\|u_j\|_{L^2}^2 = \sum_{k \in \mathbb{Z}} |\hat{u}_{j,k}|^2 = 2 \sum_{\ell=1}^{\infty} \frac{16}{2\pi(2\ell-1)^{2j+2}} = \frac{16}{\pi} \sum_{\ell=1}^{\infty} \frac{1}{(2\ell-1)^{2j+2}}$$

In order to compute this infinite series, note that

$$\zeta(s) := \sum_{n=1}^{\infty} \frac{1}{n^s},$$

is the Riemann Zeta function, and in particular for positive $s > 1$:

$$\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s} = \sum_{\ell=1}^{\infty} \frac{1}{(2\ell)^s} + \sum_{\ell=1}^{\infty} \frac{1}{(2\ell-1)^s} = \frac{1}{2^s} \zeta(s) + \sum_{\ell=1}^{\infty} \frac{1}{(2\ell-1)^s},$$

so that,

$$\sum_{\ell=1}^{\infty} \frac{1}{(2\ell-1)^s} = \left[1 - \frac{1}{2^s} \right] \zeta(s).$$

Hence, the L^2 norms we wish to compute are,

$$\|u_j\|_{L^2}^2 = \frac{16}{\pi} \sum_{\ell=1}^{\infty} \frac{1}{(2\ell-1)^{2j+2}} = \frac{16}{\pi} \left[1 - \frac{1}{2^{2j+2}} \right] \zeta(2j+2).$$

The Riemann Zeta function in particular has the following known analytical values for even integers:

$$\zeta(2j+2) = \frac{(2\pi)^{2j+2} |B_{2j+2}|}{2(2j+2)!},$$

where B_n are the Bernoulli numbers, with the following relevant values:

$$|B_n| = \begin{cases} \frac{1}{6}, & n = 2 \\ \frac{1}{30}, & n = 4 \\ \frac{1}{42}, & n = 6 \\ \frac{1}{30}, & n = 8. \end{cases}$$

Putting everything together, we conclude:

$$\|u_j\|_{L^2}^2 = \frac{16(2\pi)^{2j+1}}{(2j+2)!} \left[1 - \frac{1}{2^{2j+2}}\right] |B_{2j+2}|.$$

Hence, to numerically test L^2 convergence of $P_N u_j$ as a function of N for $j = 0, 1, 2, 3$, we plot the (square root of the) following explicit errors:

$$\begin{aligned} \|u_j - P_N u_j\|_{L^2}^2 &= \frac{16(2\pi)^{2j+1}}{(2j+2)!} \left[1 - \frac{1}{2^{2j+2}}\right] |B_{2j+2}| - \sum_{|k| \leq N} |\widehat{u}_{j,k}|^2 \\ &= \frac{16(2\pi)^{2j+1}}{(2j+2)!} \left[1 - \frac{1}{2^{2j+2}}\right] |B_{2j+2}| - \frac{16}{\pi} \sum_{\ell=1}^{\lceil N/2 \rceil} \frac{1}{(2\ell-1)^{2j+2}} \end{aligned} \quad (5)$$

We demonstrate convergence behavior using this formula in Figure 4. We see, as expected, that as j increases, the error decreases more quickly for increasing N . In particular, from the figure we also observe that the error appears to decay like,

$$\|u_j - P_N u_j\|_{L^2} \sim N^{-(j+1/2)},$$

so that at least empirically this suggests that $u_j \in H_p^{j+1/2}$. (It's actually $H_p^{j+1/2-\epsilon}$ for every $\epsilon > 0$.)

5. (Fourier interpolation)

- a.** For any $N \geq 1, k \in \mathbb{Z}$, prove that $I_N \phi_k = \phi_\ell$, where ℓ satisfying $|\ell| \leq N$ is the modular restriction of k to $[-N, N]$:

$$\ell = \ell(k) := k - (2N+1)j \in [-N, N], \quad j \in \mathbb{Z}.$$

An equivalent definition: $\ell(k) = -N + [(k+N) \pmod{2N+1}]$. The operator I_N is defined on D14-S10 as the $M = (2N+1)$ -point Fourier interpolation operator, and ϕ_k is defined on slide D13-S03(b).

- b.** If $u \in H_p^s$, prove that,

$$\|u - I_N u\|_{L^2} \lesssim N^{-s} \|u\|_{H_p^s},$$

where $a \lesssim b$ means that $a \leq Cb$ for some constant C independent of N and u .

Solution:

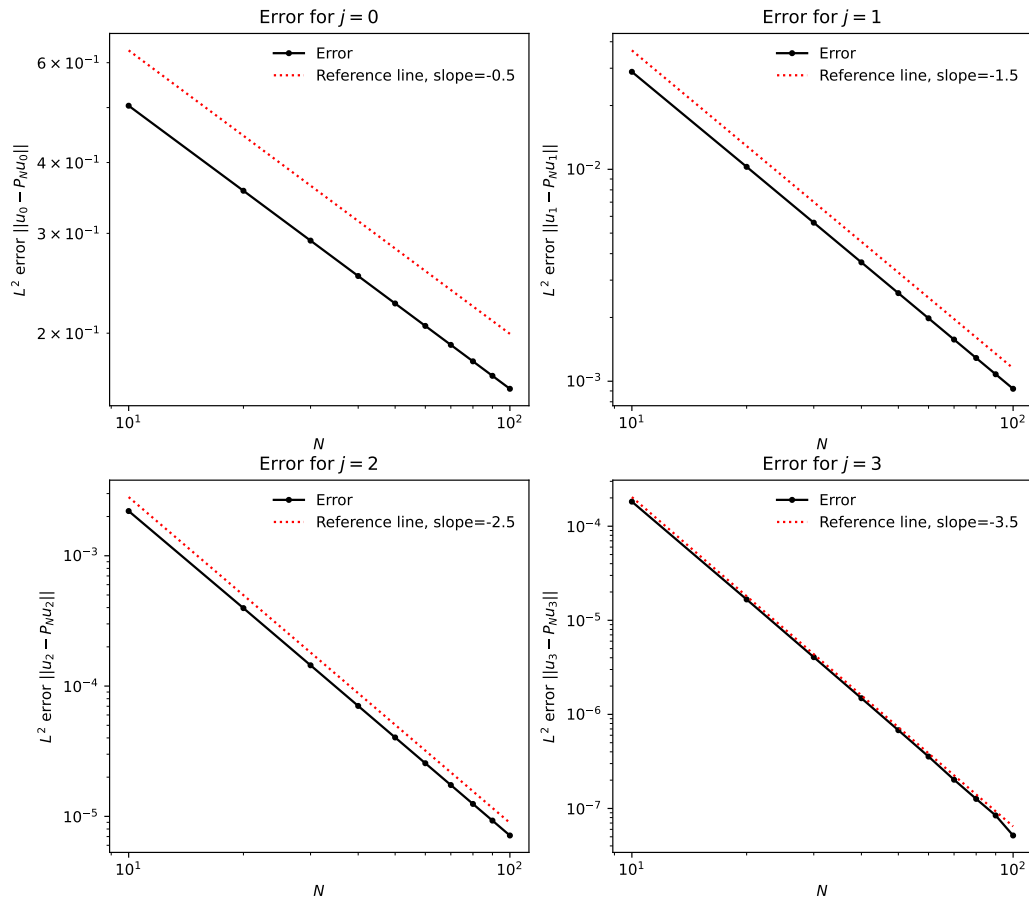


Figure 4: Problem 4: Fourier projection errors as a function of N for u_j , $j = 0, 1, 2, 3$. The errors are computed through the formula (5). Reference lines of given slopes are also plotted. One may also compute an approximate error using a truncated Parseval's sum or an interpolatory approximation to achieve essentially the same results.

- a. If $|k| \leq N$, then $k = \ell$ and ϕ_k is in the range of I_N . Since I_N is a projector, then $I_N \phi_k = \phi_k = \phi_\ell$. For $|k| > N$, then there is some $j \in \mathbb{Z}$ such that

$$\ell = k - (2N + 1)j$$

with $|\ell| \leq N$. With $x_m = 2\pi m / (2N + 1)$, $m \in [2N + 1]$ the interpolatory grid, note that,

$$\begin{aligned} \phi_k(x_m) &= \frac{1}{\sqrt{2\pi}} e^{ikx_m} = \frac{1}{\sqrt{2\pi}} e^{i\ell x_m} e^{-i(2N+1)jx_m} = \frac{1}{\sqrt{2\pi}} e^{i\ell x_j} e^{-i(2N+1)j2\pi m / (2N+1)} \\ &= \frac{1}{\sqrt{2\pi}} e^{i\ell x_j} e^{-i2\pi jm} = \frac{1}{\sqrt{2\pi}} e^{i\ell x_j} = \phi_\ell(x_m). \end{aligned}$$

Therefore, since I_N interpolates ϕ_k based only on its values at x_m and since these values equal those of ϕ_ℓ that in the range I_N , then $I_N \phi_k = \phi_\ell$. One can make this more formal: We have

$$I_N \phi_\ell = \phi_\ell = \sum_{|q| \leq N} b_q \phi_q I_N \phi_k = \sum_{|q| \leq N} c_q \phi_q$$

where we know that the coefficients \mathbf{b} are such that $b_q = 1$ with $q = \ell$ and $b_q = 0$ otherwise, but also that $\mathbf{b}, \mathbf{c} \in \mathbb{C}^{2N+1}$ are given by solutions to the quadrature/interpolation problem:

$$\begin{aligned} \mathbf{c} &= \tilde{\mathbf{V}}^* \phi_k(\mathbf{x}), \\ \mathbf{b} &= \tilde{\mathbf{V}}^* \phi_\ell(\mathbf{x}), \end{aligned}$$

where $\tilde{\mathbf{V}}$ is the scaled Fourier Vandermonde-like (DFT) matrix on slides D14-S04(b), and $\mathbf{x} = (x_m)_{m \in [2N+1]}$. Since we have already shown that $\phi_k(\mathbf{x}) = \phi_\ell(\mathbf{x})$, then $\mathbf{c} = \mathbf{b}$, implying $I_N \phi_k = I_N \phi_\ell = \phi_\ell$.

- b. Let $u = \sum_{k \in \mathbb{Z}} \hat{u}_k \phi_k$. Since the interpolation operator is linear, we have,

$$I_N u = \sum_{k \in \mathbb{Z}} \hat{u}_k I_N \phi_k = \sum_{k \in \mathbb{Z}} \hat{u}_k \phi_{\ell(k)} = \sum_{|\ell| \leq N} \left(\sum_{j \in \mathbb{Z}} \hat{u}_{\ell + j(2N+1)} \right) \phi_\ell$$

Therefore,

$$I_N u = \sum_{|k| \leq N} (\hat{u}_k + a_k) \phi_k, \quad a_k = \sum_{j \in \mathbb{Z} \setminus \{0\}} \hat{u}_{k+j(2N+1)}$$

Then the error between u and its interpolatory approximation, using the Pythagorean theorem, is,

$$\begin{aligned} \|u - I_N u\|_{L^2}^2 &\leq \|u - P_N u\|_{L^2}^2 + \|P_N u - I_N u\|_{L^2}^2 \\ &\stackrel{(*)}{\leq} N^{-2s} \|u\|_{H_p^s}^2 + \left\| \sum_{|k| \leq N} (\hat{u}_k - (\hat{u}_k + a_k)) \phi_k \right\|^2 \\ &= N^{-2s} \|u\|_{H_p^s}^2 + \sum_{|k| \leq N} |a_k|^2 \end{aligned}$$

where (*) uses the result from problem 3(a). We have shown in problem 3(a) that so long as $u \in H_p^s$, then

$$\hat{u}_k = \frac{\hat{u}_k^{(s)}}{(ik)^s},$$

where $\hat{u}_k^{(s)}$ are the Fourier expansion coefficients for $u^{(s)} \in L^2$. Using this formula for \hat{u}_k , we can bound the a_k coefficients; for example, if $k \leq 0$:

$$\begin{aligned} |a_k|^2 &= \left| \sum_{j \in \mathbb{Z} \setminus \{0\}} \hat{u}_{k+j(2N+1)} \right|^2 = \left| \sum_{j \in \mathbb{Z} \setminus \{0\}} \frac{\hat{u}_{k+j(2N+1)}^{(s)}}{(k+j(2N+1))^s} \right|^2 \\ &\stackrel{(*)}{\leq} \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \right) \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \frac{1}{|k+j(2N+1)|^{2s}} \right) \\ &\stackrel{k \leq 0}{\leq} \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \right) \left(2 \sum_{j \in \mathbb{N}} \frac{1}{(k+j(2N+1))^{2s}} \right) \\ &\stackrel{k \geq -N}{\leq} \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \right) \left(2 \sum_{j \in \mathbb{N}} \frac{1}{(jN)^{2s}} \right) \\ &= \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \right) \left(2N^{-2s} \sum_{j \in \mathbb{N}} \frac{1}{j^{2s}} \right) \\ &= 2N^{-2s} \zeta(2s) \left(\sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \right), \end{aligned}$$

where (*) is the Cauchy-Schwarz inequality (for sequences in $\ell^2(\mathbb{N})$), and $\zeta(\cdot)$ is the Riemann zeta function, which for real inputs $s > 1$ is defined as,

$$\zeta(s) := \sum_{n \in \mathbb{N}} \frac{1}{n^s}.$$

Note in particular that $\zeta(2s)$ is a finite number for any real $s \geq 1$. An analogous computation (resulting in the same bound) can be accomplished if $k > 0$. Therefore we have,

$$\begin{aligned} \|u - I_N u\|_{L^2}^2 &\leq N^{-2s} \|u\|_{H_p^s}^2 + \sum_{|k| \leq N} |a_k|^2 \\ &\leq N^{-2s} \|u\|_{H_p^s}^2 + 2\zeta(2s) N^{-2s} \sum_{|k| \leq N} \sum_{j \in \mathbb{Z} \setminus \{0\}} \left| \hat{u}_{k+j(2N+1)}^{(s)} \right|^2 \\ &= N^{-2s} \|u\|_{H_p^s}^2 + 2\zeta(2s) N^{-2s} \sum_{|k| \leq N} \|u^{(s)}\|_{L^2}^2 \\ &\leq N^{-2s} \|u\|_{H_p^s}^2 + 2\zeta(2s) N^{-2s} \|u^{(s)}\|_{L^2}^2 \\ &\leq N^{-2s} \|u\|_{H_p^s}^2 + 2\zeta(2s) N^{-2s} \|u\|_{H_p^s}^2 \\ &= (1 + 2\zeta(2s)) N^{-2s} \|u\|_{H_p^s}^2, \end{aligned}$$

which proves the result.