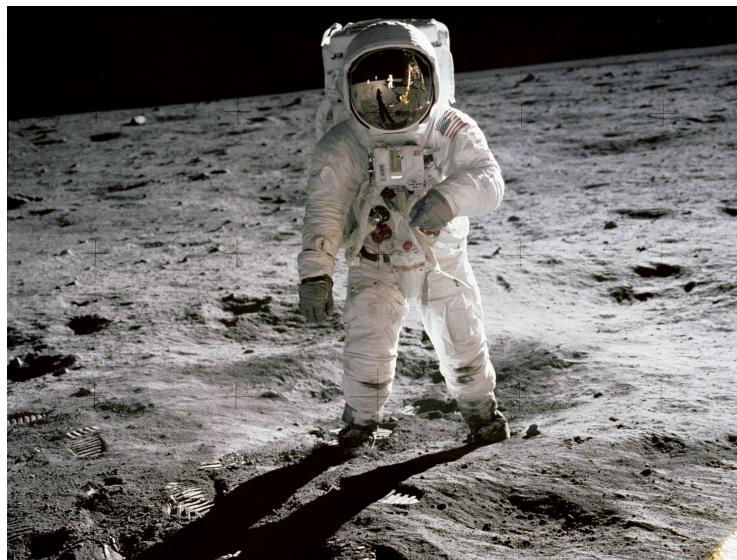


# Labs for Foundations of Applied Mathematics

Volume 4  
Modeling with Dynamics and Control

Jeffrey Humpherys & Tyler J. Jarvis, managing editors





# List of Contributors

B. Barker <i>Brigham Young University</i>	R. Dorff <i>Brigham Young University</i>
E. Evans <i>Brigham Young University</i>	B. Ehler <i>Brigham Young University</i>
R. Evans <i>Brigham Young University</i>	M. Fabiano <i>Brigham Young University</i>
J. Grout <i>Drake University</i>	K. Finlinson <i>Brigham Young University</i>
J. Humpherys <i>Brigham Young University</i>	J. Fisher <i>Brigham Young University</i>
T. Jarvis <i>Brigham Young University</i>	R. Flores <i>Brigham Young University</i>
J. Whitehead <i>Brigham Young University</i>	R. Fowers <i>Brigham Young University</i>
J. Adams <i>Brigham Young University</i>	A. Frandsen <i>Brigham Young University</i>
K. Baldwin <i>Brigham Young University</i>	R. Fuhriman <i>Brigham Young University</i>
J. Bejarano <i>Brigham Young University</i>	T. Gledhill <i>Brigham Young University</i>
A. Berry <i>Brigham Young University</i>	S. Giddens <i>Brigham Young University</i>
Z. Boyd <i>Brigham Young University</i>	C. Gigena <i>Brigham Young University</i>
M. Brown <i>Brigham Young University</i>	M. Graham <i>Brigham Young University</i>
A. Carr <i>Brigham Young University</i>	F. Glines <i>Brigham Young University</i>
C. Carter <i>Brigham Young University</i>	C. Glover <i>Brigham Young University</i>
T. Christensen <i>Brigham Young University</i>	M. Goodwin <i>Brigham Young University</i>
M. Cook <i>Brigham Young University</i>	R. Grout <i>Brigham Young University</i>

- D. Grundvig  
*Brigham Young University*
- S. Halverson  
*Brigham Young University*
- E. Hannesson  
*Brigham Young University*
- K. Harmer  
*Brigham Young University*
- J. Henderson  
*Brigham Young University*
- J. Hendricks  
*Brigham Young University*
- A. Henriksen  
*Brigham Young University*
- I. Henriksen  
*Brigham Young University*
- C. Hettinger  
*Brigham Young University*
- S. Horst  
*Brigham Young University*
- R. Howell  
*Brigham Young University*
- E. Ibarra-Campos  
*Brigham Young University*
- J. Larsen  
*Brigham Young University*
- K. Jacobson  
*Brigham Young University*
- R. Jenkins  
*Brigham Young University*
- J. Leete  
*Brigham Young University*
- Q. Leishman  
*Brigham Young University*
- J. Lytle  
*Brigham Young University*
- E. Manner  
*Brigham Young University*
- M. Matsushita  
*Brigham Young University*
- R. McMurray  
*Brigham Young University*
- S. McQuarrie  
*Brigham Young University*
- D. Miller  
*Brigham Young University*
- J. Morrise  
*Brigham Young University*
- M. Morrise  
*Brigham Young University*
- A. Morrow  
*Brigham Young University*
- R. Murray  
*Brigham Young University*
- J. Nelson  
*Brigham Young University*
- C. Noorda  
*Brigham Young University*
- A. Oldroyd  
*Brigham Young University*
- A. Oveson  
*Brigham Young University*
- E. Parkinson  
*Brigham Young University*
- M. Probst  
*Brigham Young University*
- M. Proudfoot  
*Brigham Young University*
- D. Reber  
*Brigham Young University*
- H. Ringer  
*Brigham Young University*
- C. Robertson  
*Brigham Young University*
- M. Russell  
*Brigham Young University*
- R. Sandberg  
*Brigham Young University*
- C. Sawyer  
*Brigham Young University*
- D. Smith  
*Brigham Young University*
- J. Smith  
*Brigham Young University*
- P. Smith  
*Brigham Young University*
- M. Stauffer  
*Brigham Young University*

E. Steadman

*Brigham Young University*

J. Stewart

*Brigham Young University*

S. Suggs

*Brigham Young University*

A. Tate

*Brigham Young University*

T. Thompson

*Brigham Young University*

M. Victors

*Brigham Young University*

E. Walker

*Brigham Young University*

J. Webb

*Brigham Young University*

R. Webb

*Brigham Young University*

J. West

*Brigham Young University*

R. Wonnacott

*Brigham Young University*

A. Zaitzeff

*Brigham Young University*

This project is funded in part by the National Science Foundation, grant no. TUES Phase II  
DUE-1323785.



# Preface

This lab manual is designed to accompany the textbook *Foundations of Applied Mathematics Volume 4: Modeling with Dynamics and Control* by Humpherys, Jarvis and Whitehead. The labs focus on numerical methods for solving ordinary and partial differential equations, including applications to optimal control problems. The reader should be familiar with Python [VD10] and its NumPy [Oli06, ADH<sup>+</sup>01, Oli07] and Matplotlib [Hun07] packages before attempting these labs. See the Python Essentials manual for introductions to these topics.

©This work is licensed under the Creative Commons Attribution 3.0 United States License. You may copy, distribute, and display this copyrighted work only if you give credit to Dr. J. Humpherys. All derivative works must include an attribution to Dr. J. Humpherys as the owner of this work as well as the web address to

<https://github.com/Foundations-of-Applied-Mathematics/Labs>  
as the original source of this work.

To view a copy of the Creative Commons Attribution 3.0 License, visit

<http://creativecommons.org/licenses/by/3.0/us/>  
or send a letter to Creative Commons, 171 Second Street, Suite 300, San Francisco, California, 94105, USA.





# Contents

Preface	v
<b>I    Labs</b>	<b>1</b>
1    Introduction to Matplotlib: 3D Plotting and Animations	3
2    Intro to IVP and BVP	11
3    Modelling the spread of an epidemic: SIR models	23
4    Numerical Methods for Initial Value Problems; Harmonic Oscillators	35
5    Weight change and Predator-Prey Models	45
6    Lorenz Equations	53
7    Bifurcations	59
8    The Finite Difference Method	69
9    Wave Phenomena	77
10   Heat Flow	87
11   Anisotropic Diffusion	95
12   The Finite Element method	103
13   Poisson's equation	111
14   Method of Mean Weighted Residuals	119
15   A Pseudospectral method for periodic functions	125
16   Inverse Problems	131
17   The Shooting Method for Boundary Value Problems	137

<b>18</b>	<b>Total Variation and Image Processing</b>	<b>147</b>
<b>19</b>	<b>Transit time crossing a river</b>	<b>155</b>
<b>20</b>	<b>HIV Treatment Using Optimal Control</b>	<b>161</b>
<b>21</b>	<b>Solitons</b>	<b>169</b>
<b>22</b>	<b>Obstacle Avoidance</b>	<b>175</b>
<b>23</b>	<b>The Inverted Pendulum</b>	<b>183</b>
<b>II</b>	<b>Appendices</b>	<b>191</b>
<b>A</b>	<b>Getting Started</b>	<b>193</b>
<b>B</b>	<b>Installing and Managing Python</b>	<b>201</b>
<b>C</b>	<b>NumPy Visual Guide</b>	<b>207</b>
<b>D</b>	<b>Matplotlib Customization</b>	<b>211</b>
	<b>Bibliography</b>	<b>227</b>

# Part I

## Labs



# 1

# Introduction to Matplotlib: 3D Plotting and Animations

**Lab Objective:** Animations and 3D plots are useful in visualizing solutions to ODEs and PDEs found in many dynamics and control problems. In this lab we explore the functionality contained in the 3D plotting and animation libraries in Matplotlib.

## Introduction

Matplotlib is a Python library that contains tools for creating plots in multiple dimensions. The library contains important classes that are needed to create plots. The most important objects to understand in this lab are figure objects, axes objects, and line objects. These three objects are created using the following code.

```
>>> import matplotlib.pyplot as plt
>>> fig = plt.figure()                      # Create figure object.
>>> ax = fig.add_subplot(111)                 # Create axes object.
>>> line2d, = plt.plot([],[])                # Create empty 2D Line object
>>> line3d, = plt.plot([],[],[])             # Create empty 3D line object
```

Recall that `plt.figure()` creates a `matplotlib.figure.Figure` object, which is the window that is displayed when `plt.show()` is called. 3D plotting and animation both require explicitly defining the `Figure` object, as shown above. This allows for the object to be updated and modified, as will be explained later in the lab.

`Figure` objects contain `matplotlib.axes._subplots.AxesSubplot` objects, called `axes`. `Axes` are spaces to plot on, and are created by the `add_subplot()` method of a `Figure` object. Figures can have multiple axes.

Calling `plt.plot()` returns a list of line objects. For example, supposing `x1`, `y1`, `x2`, and `y2` are arrays containing data for two separate curves, then calling `plt.plot(x1, y1, x2, y2)` will return a list with two elements. Each element of the list is a `matplotlib.lines.Line2D` object. If the axes is three-dimensional, then the returned list will contain `matplotlib.lines.Line3D` objects. Because this function call returns a list, if only one line is plotted, adding a trailing comma to the variable name will assign the name to the first element of the returned list. You can alternatively reference the zero index of the returned list, but using a trailing comma is standard.

## Animation Basics

The animation library in Matplotlib contains a class called `FuncAnimation`. We will use this class throughout this lab. `FuncAnimation` requires a user-defined update function that controls the plot for each frame of the animation. This grants the user wide flexibility and control of the resulting animation. The following steps describe the process of creating a simple animated plot using the `FuncAnimation` class.

1. Compute all data to be plotted.
2. Explicitly define figure object.
3. Define line objects to be altered dynamically.
4. Create function to update line objects.
5. Create `FuncAnimation` object.
6. Display using `plt.show()`.

These steps will be explained by way of an example. The arrays `x` and `y` contain data giving the location of a particle moving in the plane. To visualize this motion, one could animate the particle as well as display the trajectory that the particle has traveled. For this animation, two separate `Line2D` objects must be created on an axes object. The first, `particle` will be for the position of the particle itself, and the second, `traj` will be for the trajectory that the particle has traveled. Note that these objects are created with empty lists of data. The update function will be used to dynamically set the data to be plotted in these line objects.

```
>>> import matplotlib.animation as animation
>>> import numpy as np
>>> t = np.linspace(0,2*np.pi,100)
>>> x = np.sin(t)
>>> y = t**2
>>> fig = plt.figure()
>>> ax = fig.add_subplot(111)
>>> ax.set_xlim((-1.1,1.1))
>>> ax.set_ylim((0,40))
>>> particle, = plt.plot([],[], marker='o', color='r')
>>> traj, = plt.plot([],[], color='r', alpha=0.5)
```

The update function must be defined a specific way in order to interact properly with the `matplotlib.animation.FuncAnimation` object. The update function must accept the current frame index as its first input parameter and it must return a list or tuple of line objects. The current frame index is used to access the data to be plotted in the current frame. Both 2D and 3D line objects have the built-in method `.set_data()`. This function takes in two one-dimensional arrays representing `x` and `y` values to plot. This allows a single line object to display different data for each frame. Inside the update function, `.set_data()` is called on the line objects with the relevant data as inputs.

```
>>> def update(i):
>>>     particle.set_data(x[i],y[i])
>>>     traj.set_data(x[:i+1],y[:i+1])
>>>     return particle,traj
```

Next, the `FuncAnimation` object is created. The argument `frames` specifies the iterable representing the frame indices. If `frames` is an integer, it is treated as the iterable `range(frames)`. After the `FuncAnimation` object is created, `plt.show()` displays the animation.

```
>>> ani = FuncAnimation(fig, update, frames=range(100), interval=25)
>>> plt.show()
```

The following table shows more parameters that can be passed into `FuncAnimation`.

Parameter	Description
<code>fargs</code> (tuple)	Additional arguments to pass update function
<code>interval</code> (float)	Delay between frames in milliseconds
<code>repeat</code> (bool)	Determines whether animation repeats (Default True)
<code>blit</code> (bool)	Determines whether blitting is used (Default False)

### Note

When using `FuncAnimation`, it is essential that a reference is kept to the instance of the class. The animation is advanced by a timer and if a reference is not held for the object, Python will automatically garbage collect and the animation will stop.

### Saving Animations

The simplest way to save an animation is to encode it to a `.mp4` file, which will allow you to display the video inline inside a Jupyter Notebook, or view it using any video player supporting the chosen filetype.

Unfortunately, Matplotlib does not come with a built-in video encoder. The `matplotlib.animation` module supports several third-party encoders. FFmpeg is a lightweight solution which can be obtained from <https://www.ffmpeg.org/download.html>. When available, FFmpeg is generally chosen as the default, but you may need to specify to Matplotlib to use it:

```
animation.writer = animation.writers['ffmpeg']
```

To prevent the animation from displaying while it is being rendered as video, use `plt.ioff()`. This turns off matplotlib's interactive mode until `plt.ion()` is called. After creating the animation object, use its `.save()` method with the desired filename to render and save the video. The following code is given for reference:

```
plt.ioff()      # Turn off interactive mode to hide rendering animations

# Code to create figure, axes, and update function goes here
# ...
ani = animation.FuncAnimation(fig, update, frames, interval)
ani.save('my_animation.mp4')
```

To display the `.mp4` video in a Jupyter Notebook, place the following HTML code in a separate markdown cell:

```
<video src="my_animation.mp4" controls>
```

## Embedding Animations

While saving animations to a file has the advantage that the animation will persist if the notebook is closed and reopened, it tends to be much slower than directly embedding the animation in the notebook. Directly embedding can thus be useful in the process of creating an animation by allowing faster experimentation.

After creating an animation, calling `plt.show()` will attempt to embed it; however, some systems may struggle to display an animation in this way. When this is the case, it may be easier to embed the animation using the HTML5 API. Jupyter notebooks use HTML to display their contents, so we can leverage this and use HTML5's video capabilities to insert video directly into a notebook.

To embed the video directly into a notebook using HTML5 you must use the `IPython.display` module. This module will be able to interpret an encoded HTML5 video, which `matplotlib.animation` can create. This method tends to be much more simple than rendering the animation to an `.mp4` file and then embedding that file into a notebook, as it does not require an outside encoder, and it tends to be encounter fewer bugs than using `plt.show()`. However, the animation generally does not persist if the notebook is closed and reopened. Here is a snippet you may reference to embed an animation using `IPython.display`

```
# required import statements
from IPython.display import HTML
import matplotlib.pyplot as plt
from matplotlib import animation

# disable interactive mode
plt.ioff()
'''

Here we would insert whatever code needed to create the animation
such as instantiating the fig object and defining the update function
'''

# create animation
ani = animation.FuncAnimation(fig, update, frames, interval)
# render as html5 and embed
HTML(ani.to_html5_video())
```

### Achtung!

Note that animations that are embedded in the notebook using `plt.show()` or `HTML()` **do not** persist if the notebook is closed and reopened. While this method can be useful for more quickly testing animations, **do not** use this method for embedding the final animations of your finished lab, as the grader will not be able to view your animations. Rather, save the animation to an `.mp4` file and embed the created file. When pushing your lab, be sure to also add and push the video files you create.

**Problem 1.** Use the FuncAnimation class to animate the function  $y = \sin(x + 0.1t)$  where  $x \in [0, 2\pi]$ , and  $t$  ranges from 0 to 100 seconds. Save your animation to a file and embed that into the notebook.

## 3D Plotting Introduction

3D plotting is very similar to 2D plotting. The main difference is that a set of 3D axes must be created within the figure object. A 3D axes object is created using the additional keyword argument `projection='3d'`:

```
>>> # Create figure object.
>>> fig = plt.figure()
>>>
>>> # Create 3D axis object using add_subplot().
>>> ax = fig.add_subplot(111, projection='3d')
```

## 3D Static Plotting

When the axes object is explicitly defined, plots are generated by calling the chosen plot function (such as `ax.plot()` on the axes object. Additional information on the use of axes objects can be found here: [https://matplotlib.org/api/axes\\_api.html](https://matplotlib.org/api/axes_api.html).

**Problem 2.** The orbits for Mercury, Venus, Earth, and Mars are stored in the file `orbits.npz`. The file contains four NumPy arrays: `mercury`, `venus`, `earth`, and `mars`. The first column of each array contains the x-coordinates, the second column contains the y-coordinates, and the third column contains the z-coordinates of each planet, all relative to the Sun, and expressed in AU (astronomical units, the average distance between Earth and the Sun, approximately 150 million kilometers).

Use `np.load('orbits.npz')` to load the data for the four planets' orbits. Create a 3D plot of the orbits, and compare your results with Figure 1.1.

## 3D Animations

The key difference between 2D and 3D animations is that the `.set_data()` method does not support setting the `z` values. Instead, set the `x` and `y` values with `.set_data()` as before, and then set the `z` values with `.set_3d_properties()`. The `.set_3d_properties()` function call is also made inside the update function.

Animation in 3D requires more careful consideration than in the 2D case. When `matplotlib` displays a 3D plot, it does so in an interactive figure that allows the user to change the camera angle and position. Since 3D rendering is more computationally expensive than 2D rendering, interactive views of 3D animations often have poor framerates and choppy rendering. This is what calling `plt.plot()` attempts to do; instead, it is much better to either render the animation to a file and then embed the file, or to use the HTML5 API to embed it, as discussed above.

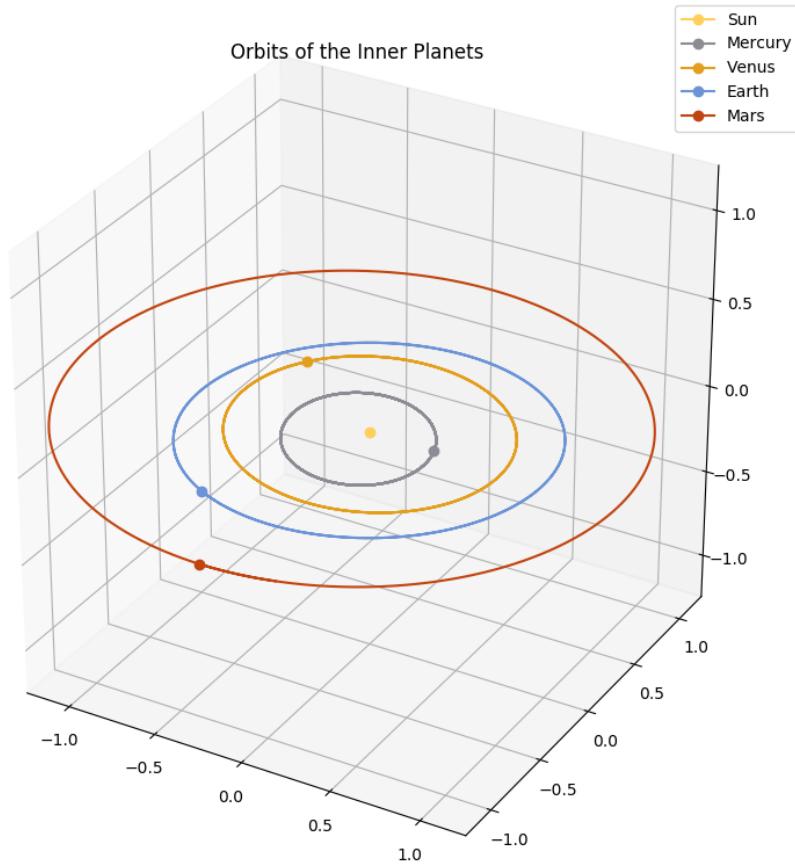


Figure 1.1: The solution to Problem 2.

**Problem 3.** Each row of the arrays in `orbits.npy` gives the position of the planets at evenly spaced time points. The arrays correspond to 1400 points in time over a 700 day period (beginning on 2018-5-30). Create a 3D animation of the planet orbits. Display lines for the trajectories of the orbits and points for the current positions of the planets at each point in time. Your `update()` function will need to return a list of `Line3D` objects, one for each orbit trajectory and one for each planet position marker. Embed your animated plot.

## Surface Plotting

3D surface plotting is very similar to regular 3D plotting discussed earlier. The difference with surface plots is that they require first creating a meshgrid for X and Y. Meshgrids are created using the NumPy command `np.meshgrid(x, y)` where `x` and `y` are 1D arrays representing the x and y coordinates of the grid. This function creates 2D arrays `X` and `Y` that combined give cartesian coordinates for every point made from the `x` and `y` arrays.

Once a meshgrid is defined, a surface plot is generated by calling `ax.plot_surface(X, Y, Z)`, where `Z` is a 2D array of height values that is the same shape as `X` and `Y`.

**Problem 4.** Make a surface plot of the bivariate normal density function given by:

$$f(\mathbf{x}) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

where  $\mathbf{x} = [x, y]^T$ ,  $\boldsymbol{\mu} = [0, 0]^T$  is the mean vector, and

$$\Sigma = \begin{bmatrix} 1 & 3/5 \\ 3/5 & 2 \end{bmatrix}$$

is the covariance matrix. Compare your results with Figure 1.2.

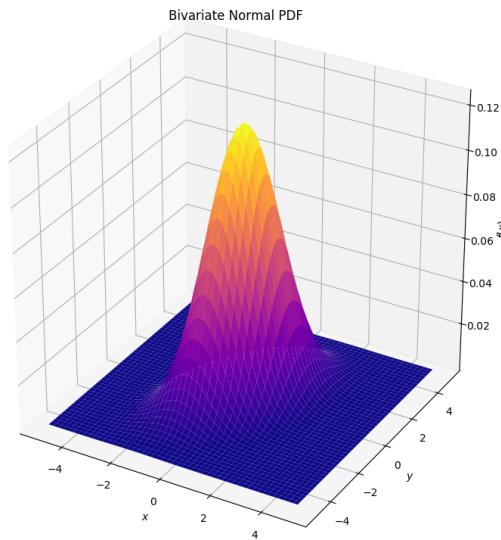


Figure 1.2: The solution to Problem 4.

## Surface Animations

Animating a 3D surface is slightly different from animating a parametric curve in 3D. The object created by `.plot_surface()` does not have a `.set_data()` method. Instead, use `ax.clear()` to empty the axes at each frame, followed by a new call to `ax.plot_surface()`. Note that the axes limits must be reset after `ax.clear()` is called.

**Problem 5.** Use the data in `vibration.npz` to produce a surface animation of the solution to the wave equation for an elastic rectangular membrane. The file contains three NumPy arrays: `X`, `Y`, `Z`. `X` and `Y` are meshgrids of shape `(300, 200)` corresponding to 300 points in the `y`-direction and 200 points in the `x`-direction, giving a  $2 \times 3$  rectangle with one corner at the origin. `Z` is of shape `(150, 300, 200)`, giving the height of the vibrating membrane at each  $(x, y)$  point for 150 values of time. In the language of partial differential equations, this is the solution to the following initial/boundary value problem:

$$\begin{aligned} u_{tt} &= 6^2(u_{xx} + u_{yy}) \\ (x, y) &\in [0, 2] \times [0, 3], t \in [0, 5] \\ u(t, 0, y) &= u(t, 2, y) = u(t, x, 0) = u(t, x, 3) = 0 \\ u(0, x, y) &= xy(2 - x)(3 - y) \end{aligned}$$

# 2

## Intro to IVP and BVP

### Initial Value Problems

An initial value problem is a differential equation with a set of constraints at the initial point. An IVP may look something like this

$$\begin{aligned}y'' + y' + y &= f(t) \\y(a) &= \alpha \\y'(a) &= \beta \\t \in [a, b].\end{aligned}$$

This problem gives a differential equation with initial conditions for  $y$  and  $y'$ .

Formulating and solving initial value problems is an important tool when solving many types of problems. One simple example of an IVP would be a differential equation modeling the path of a ball thrown in the air where the initial position ( $y(a)$ ) and velocity ( $y'(a)$ ) are known. These problems can be tricky to solve by hand. Luckily, SciPy has great tools that help us solve initial value problems for most systems of first order ODEs. We will be using `solve_ivp` from `scipy.integrate`.

Consider the following example

$$y'' + 3y = \sin(t), \quad y(0) = -\pi/2, \quad y'(0) = \pi, \quad t \in [0, 5]$$

We begin by changing this second order ODE into a first order ODE system.

Let  $y_1 = y$  and  $y_2 = y'$  so that

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}' = \begin{bmatrix} y_2 \\ \sin t - 3y_1 \end{bmatrix}.$$

This formulation allows us to use `solve_ivp`. We need three code elements in order to use `solve_ivp`:

1. The ODE function:

This function defines the right-hand side of the ODE system, and returns an array containing its values for our first order system of ODEs.

2. The time domain:

This is a tuple giving the interval of integration.

3. The initial conditions:

This is an array containing the initial conditions of each coordinate of the ODE. In our example, these are the value of the "zeroeth" derivative, followed by the first derivative, and so on if there are higher order derivatives.

The following code sets up and solves the IVP in the above example:

```
from scipy.integrate import solve_ivp
import numpy as np

# element 1: the ODE function
def ode(t, y):
    '''defines the ode system'''
    return np.array([y[1],np.sin(t)-3*y[0]])

# element 2: the time domain
t_span = (0,5)

# element 3: the initial conditions
y0 = np.array([-np.pi /2, np.pi])

# solve the system
# max_step is an optional parameter that controls maximum step size and
# a smaller value will result in a smoother graph
sol = solve_ivp(ode, t_span, y0, max_step=0.1)

# as an alternative, the parameter t_eval can be used to evaluate the function
# at specific points; this can also be used to get a smooth graph
sol = solve_ivp(ode, t_span, y0, t_eval=np.linspace(-np.pi/2, np.pi, 150))
```

Then we can plot the solution with the following code:

```
from matplotlib import pyplot as plt

plt.plot(sol.t,sol.y[0])
plt.xlabel('t')
plt.ylabel('y(t)')
plt.show()
```

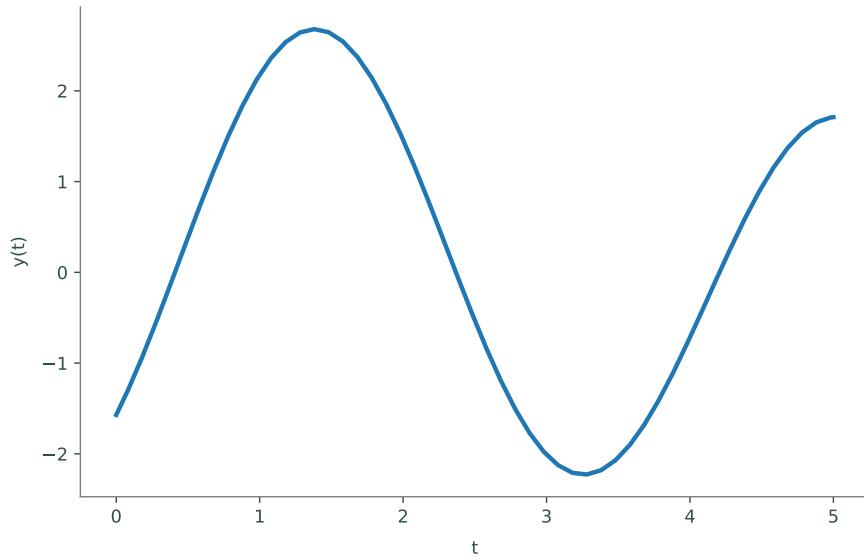


Figure 2.1: The solution to the above example

**Problem 1.** Use `solve_ivp` to solve for  $y$  in the equation  $y'' - y = \sin(t)$  with initial conditions  $y(0) = -\frac{1}{2}$ ,  $y'(0) = 0$  and plot your solution on the interval  $[0, 5]$ . Compare this to the analytic solution  $y = -\frac{1}{2}(e^{-t} + \sin(t))$ .

Note: Using `max_step = 0.1` will give you the smoother graph seen in figure 2.2.

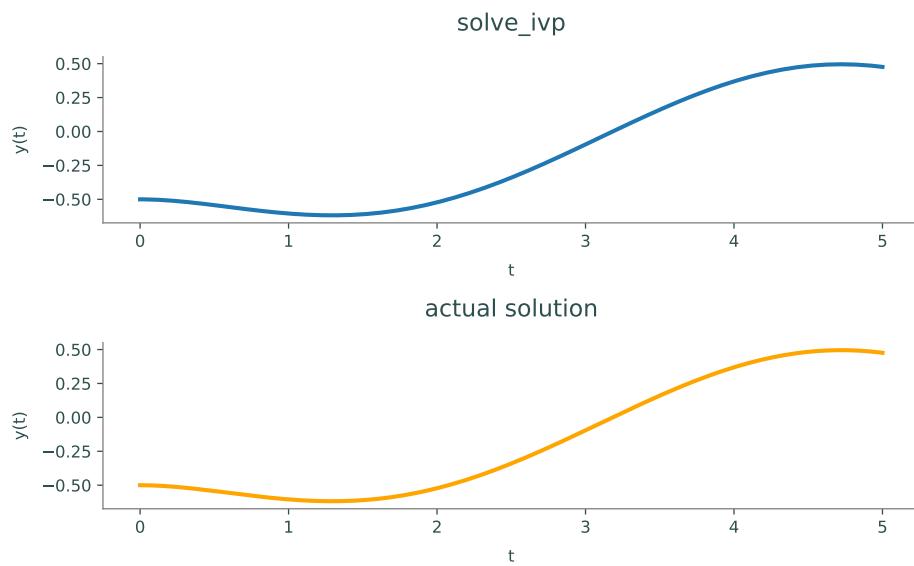


Figure 2.2: The solution to Problem 1

## Boundary Value Problems

A boundary value problem is a differential equation with a set of constraints. It is similar to initial value problems, but may give end constraints as well as initial constraints. A boundary value problem may look something like this

$$\begin{aligned}y'' + y' + y &= f(t) \\y(a) &= \alpha \\y(b) &= \beta \\t \in [a, b],\end{aligned}$$

where we have both right and left hand boundary conditions on  $y$ . One simple example of an IVP would be a differential equation modeling the path of a ball thrown in the air where the initial position ( $y(a)$ ) and final position ( $y(b)$ ).

SciPy has great tools that help us solve boundary value problems. We will be using `solve_bvp` from `scipy.integrate`. Consider the following example:

$$y'' + 9y = \cos(t), \quad y'(0) = 5, \quad y(\pi) = -\frac{5}{3}. \quad (2.1)$$

We begin by changing this second order ODE into a first order ODE system.

Let  $y_1 = y$  and  $y_2 = y'$  so that,

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}' = \begin{bmatrix} y_2 \\ \cos t - 9y_1 \end{bmatrix}.$$

This formulation allows us to use `solve_bvp`. It is important to notice that there are several key differences between `solve_ivp` and `solve_bvp`. We need four code elements in order to use `solve_bvp`:

1. The ODE function:

This is essentially the same function we used in `solve_ivp`.<sup>1</sup>

2. The boundary condition function:

Instead of just having a tuple containing our initial values, we now must use a function that returns an array of the residuals of the boundary conditions. We pass in 2 arrays: `ya`, representing the initial values, and `yb`, representing the final values. The  $i$ th entry of those arrays represents the boundary condition at the  $i$ th coordinate of the ODE. Returning `ya[0]-x` would indicate that we know  $y_1(a) = x$ , `ya[1]-x` would indicate that we know  $y_2(a) = x$ , `yb[0]-x` would indicate that we know  $y_1(b) = x$ , and `yb[1]-x` would indicate that we know  $y_2(b) = x$ .

3. The time domain:

Instead of a tuple giving the interval of integration, we now must pass in a `linspace` from the starting time to the ending time, containing the desired number of points (we now must choose the number). As part of its algorithm, `solve_bvp` will add additional points to the mesh to attempt to reduce the error of the approximation, so it is not generally necessary to pass in a very fine mesh. This also means that the mesh of the returned solution will generally not be the same as the one you pass in here.

---

<sup>1</sup>There is a technical difference between how the two methods call the ODE function. Unlike `solve_ivp`, `solve_bvp` calls the function on all of the time steps all at once, so `t` will be an array and `y` will be a  $(n, T)$  array where  $n$  is the dimension of the ODE and  $T$  is the number of timesteps. For most applications, this leads to no difference in how you code the ODE function, as can be seen in the examples; however, for some applications, such as piecewise ODE functions, this fact must be taken into consideration.

4. The initial guess:

As we no longer know all of the initial values, we now must make a (hopefully educated) guess. This is an array of shape ( $n$ ,  $t\_steps$ ) where  $n$  is the shape of the output of the ODE function and  $t\_steps$  is the chosen number of steps in our time domain linspace.

```
from scipy.integrate import solve_bvp
import numpy as np

# element 1: the ODE function
def ode(t,y):
    ''' define the ode system '''
    return np.array([y[1], np.cos(t) - 9*y[0]])

# element 2: the boundary condition function
def bc(ya,yb):
    ''' define the boundary conditions '''
    # ya are the initial values
    # yb are the final values
    # each entry of the return array will be set to zero
    return np.array([ya[1] - 5, yb[0] + 5/3])

# element 3: the time domain.
t_steps = 100
t = np.linspace(0,np.pi,t_steps)

# element 4: the initial guess.
y0 = np.ones((2,t_steps))

# Solve the system.
sol = solve_bvp(ode, bc, t, y0)
```

The syntax for plotting the function is also slightly different:

```
import matplotlib.pyplot as plt

# here we plot sol.x instead of sol.t
plt.plot(sol.x, sol.y[0])
plt.xlabel('t')
plt.ylabel('y(t)')
plt.show()
```

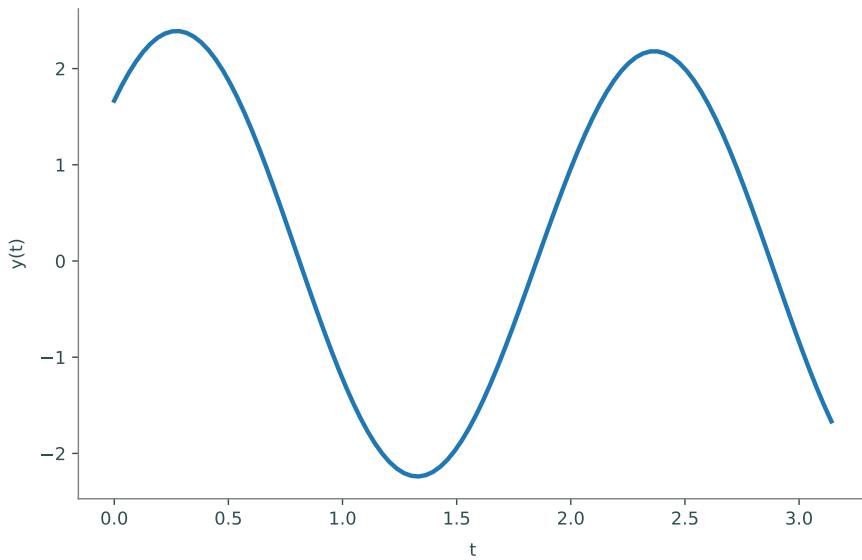


Figure 2.3: The solution to the above boundary value problem

**Problem 2.** Use `solve_bvp` to solve for  $y$  in the equation  $y'' + y' = -\frac{1}{4}e^{-t/2} + \sin(t) - \cos(t)$  with boundary conditions  $y(0) = 6$ ,  $y'(5) = -0.324705$  and plot your solution on the interval  $[0, 5]$ . Compare this to the analytic solution  $y = e^{-t/2} - \sin(t) + 5$ .

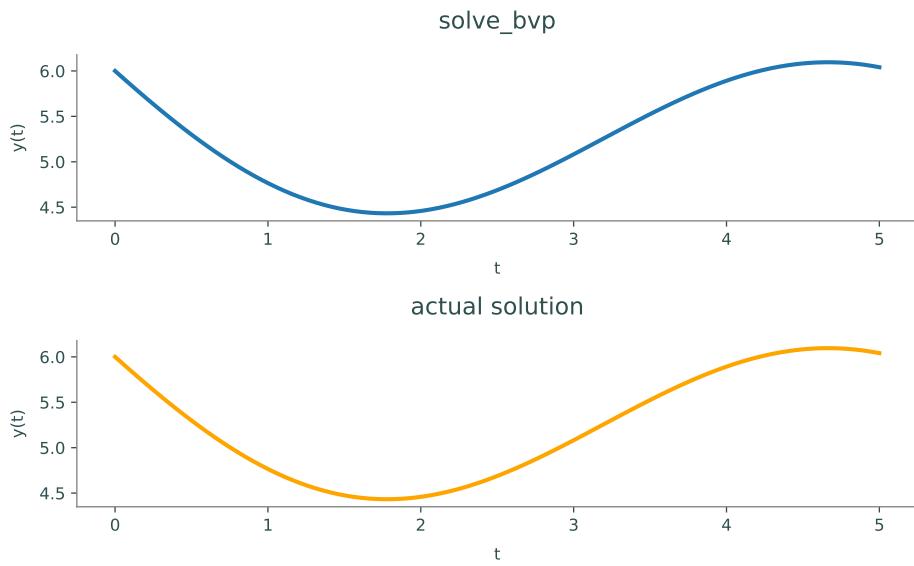


Figure 2.4: The solution to problem 2.

One other useful functionality of `solve_bvp`: `sol.sol` is a callable function, which is the estimation of the boundary value problem. You can plug in any value or numpy array (`sol.sol(np.linspace)`, `sol.sol(float)`, `sol.sol(list)`), like a normal lambda function.

## The Pitfalls of `solve_bvp`

One of the common issues with `solve_bvp` is choosing a guess for the initial value. Often, small changes in the guess can cause large changes in the final approximation. The reason for this is that the algorithm used by `solve_bvp` is essentially a version of Newton's method set up to approximate the boundary value problem, and thus can be sensitive to the initial guess. The next problem demonstrates the huge difference that can be made between a constant initial guess of 10 and a constant initial guess of 9.99

**Problem 3.** Use `solve_bvp` to solve for  $y$  in the equation  $y'' = (1 - y') * 10y$  with boundary conditions  $y(0) = -1$  and  $y(1) = \frac{3}{2}$  and plot your solution on the interval  $[0, 1]$ . Use an initial guess of 10. Compare this to the same solution using an initial guess of 9.99. For both of your initial guesses, use 50 steps in  $t$ .

The solution is found in Figure 2.5.

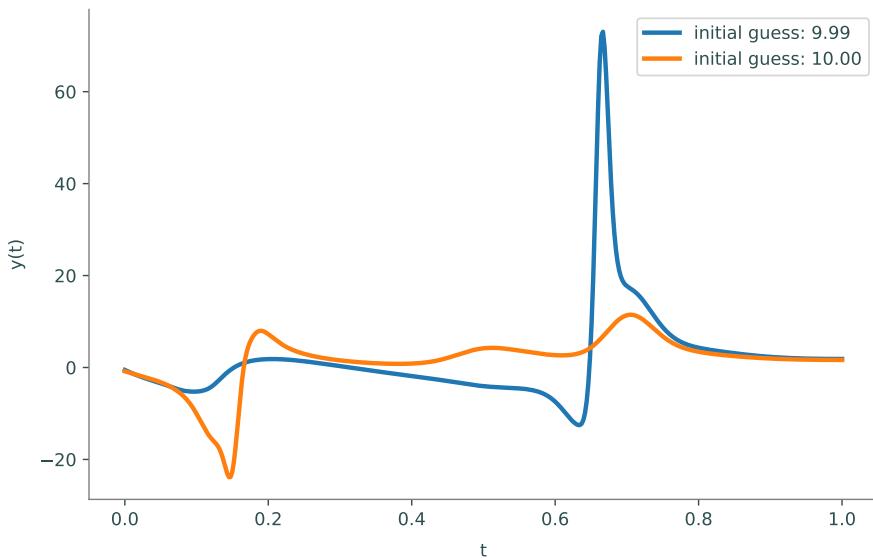


Figure 2.5: The solution to problem 3.

## Strange Attractors

In the growing field of dynamical systems, an attractor is a set of states toward which a system tends to evolve. Strange attractors are special in that they showcase complex behavior in a simple set of equations. A minute change in the initial values can cause massive differences in the outcome. The most famous of these is the Lorenz attractor, introduced by Edward Lorenz in 1963. Later on, we dedicate a full lab to the study of the Lorenz attractor, but today we focus on modeling the Four-Wing attractor. This is a system of first order ODEs defined by the following set of equations

$$\frac{dx}{dt} = ax + yz \quad (2.2)$$

$$\frac{dy}{dt} = bx + cy - xz \quad (2.3)$$

$$\frac{dz}{dt} = -z - xy \quad (2.4)$$

given some constants  $a$ ,  $b$ , and  $c$ . As we mentioned earlier, `solve_ivp` and `solve_bvp` can be used to solve and model systems of first order ODEs. We will now use `solve_ivp` to model the Four-Wing attractor.

**Problem 4.** Use `solve_ivp` to solve the Four-Wing Attractor as described in equations (1.2), (1.3), and (1.4) where  $a = 0.2$ ,  $b = 0.01$ , and  $c = -0.4$ . Try this with 3 different initial values and plot (in three dimensions) the 3 corresponding graphs.

Examples of solutions are given in Figure 2.6.

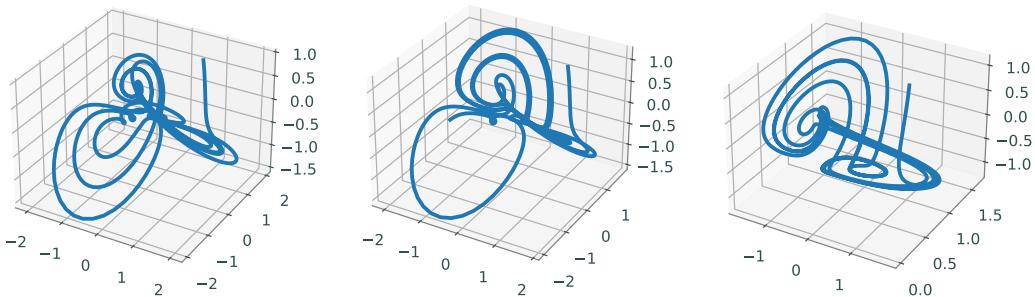


Figure 2.6: Possible solutions to problem 4.

## BVPs with Free Parameters

The function `solve_bvp` also supports solving for free parameters in a boundary value problem. The syntax is very similar to what we did above, except that our functions will have an additional argument that is a list of all of the free parameters. For example, suppose we have the following BVP with free parameter  $a$ :

$$\begin{aligned} x'(t) &= (1+t)y(t) - x(t) \\ y'(x) &= a \sin(x(t)) \\ x(0) &= 1, \quad x(1) = 0, \quad y(1) = 2 \end{aligned}$$

Note that we need one additional boundary condition for each free parameter; in this case, since we have two variables and one free parameter, we need three boundary conditions. We set up the ODE and boundary condition functions as follows:

```
# The ODE function
def ode(t, y, p):
    ''' Defines the ODE system '''
    return np.array([
        (1+t)* y[1] - y[0],
        p[0] * np.sin(y[0])
    ])

# The boundary condition function
def bcs(ya, yb, p):
    ''' Defines the boundary conditions '''
    return np.array([
        ya[0] - 1,
        yb[0] - 0,
        yb[1] - 2
    ])
```

Note that both of these functions accept an additional argument `p`, which is a list of the free parameters in the problem. In this case, we only have one parameter, so `p=[a]`. Using `solve_bvp` to get the solution is similar to before, except that we must also pass in a guess for the free parameters with the argument `p`:

```
# Guess of the solution values
t = np.linspace(0, 1, 50)
y_guess = np.ones((2,50))
p_guess = [1]

# Solve
sol = solve_bvp(ode, bcs, t, y_guess, p=p_guess)
```

The solution can be plotted as before, and the value of the free parameters for the solution can be found with `sol.p`:

```
plt.plot(sol.x, sol.y[0], label='$x(t)$')
plt.plot(sol.x, sol.y[1], label='$y(t)$')
plt.legend()
plt.xlabel('t')
plt.ylabel('y(t)')
plt.title(f'$a = {sol.p[0]:.4f}$')
plt.show()
```

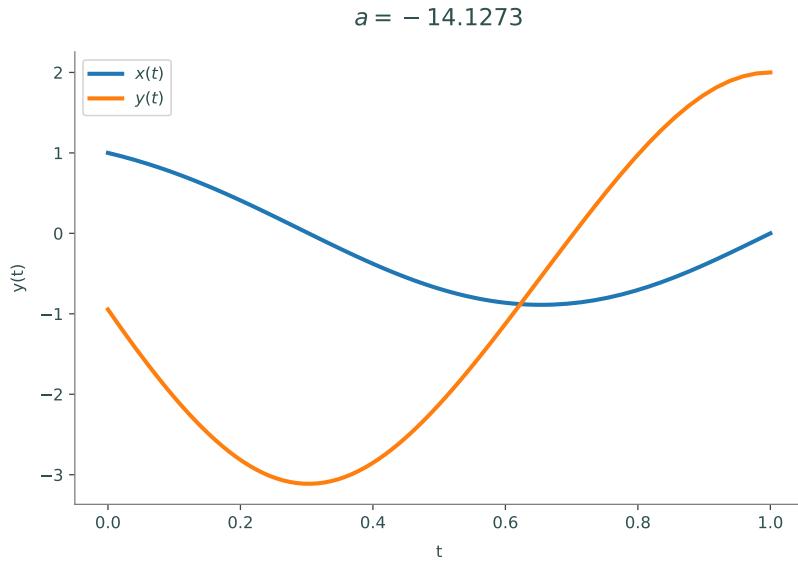


Figure 2.7: The solution to the above example

Free parameters occur in eigenvalue problems for differential operators. The general form of these is

$$\begin{aligned} Dy &= \lambda y \\ y(a) &= y(b) = 0 \end{aligned}$$

where  $D$  is some differential operator and  $\lambda$  is an unknown scalar. These differential eigenvalue problems are very analogous to the finite-dimensional case of matrix eigenvalue problems, except that instead of trying to find eigenvalues, we now try to find eigenfunctions. As in the matrix case, we are not interested in the trivial solution  $y = 0$ . Furthermore, if we have an eigenfunction  $y$ , any multiple  $ay$  is also an eigenfunction. To solve both of these issues, we will stipulate that  $y'(a) = 1$ ; this guarantees the solution is not identically zero, and also makes it unique for a given value of  $\lambda$ .

Sturm-Liouville problems are an important category of these eigenvalue problems. These have the special form

$$\begin{aligned} (py')' + qy &= \lambda ry \\ y(a) &= y(b) = 0 \end{aligned}$$

where  $p(t), q(t), r(t)$  are known functions and  $\lambda$  is an unknown scalar. Sturm-Liouville problems and their extensions are important theoretically, and their solutions have some very nice properties. They also have applications in PDEs and occur in areas such as physics and quantum mechanics.

**Problem 5.** An important problem in quantum mechanics is to find steady-state solutions of the Schrödinger equation. These functions are solutions to the time-independent Schrödinger equation. This equation is a differential equation for the wave function  $\psi$ , with one free parameter  $E$ . In one dimension, this equation is

$$-\frac{\hbar^2}{2m}\psi''(x) + U(x)\psi(x) = E\psi(x). \quad (2.5)$$

where  $U$  is a known function describing the potential energy. Note that this is in fact a Sturm-Liouville problem. If a function  $\psi$  and scalar  $E$  satisfy this equation, they describe an allowed steady state of a particle in the system. The value of  $E$  is the energy of the particle in that state, and the values of  $\psi$  are related to the probability of the particle being in any given location. For simplicity, we will let  $\frac{\hbar^2}{2m} = 1$ .<sup>a</sup>

Write a function that uses `solve_bvp` to find  $\psi$  and  $E$  that are solutions to (2.5) for the potential  $U(x) = x^2$  and with boundary conditions  $\psi(-1) = \psi(1) = 0, \psi'(-1) = 1$ . By varying your initial guess for  $E$ , use your function to find solutions for several different values of  $E$ , and plot them together.

---

<sup>a</sup>Making constants equal to one in this way is actually done quite frequently in particle physics, by choosing the units we are using appropriately.

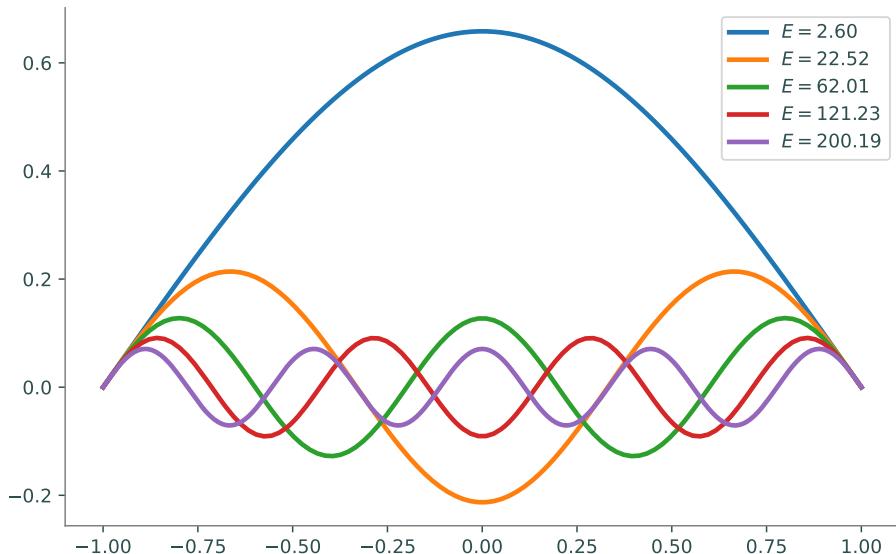


Figure 2.8: A possible solution to problem 5.



# 3

## Modelling the spread of an epidemic: SIR models

### Numerical Solvers

We often rely on numerical solvers to numerically integrate ordinary differential equations, ODEs. Because of the complexity of many ODE systems, these numerical solvers allow us to solve complex ODE systems that may not be solvable symbolically, or are high dimensional. In this lab we will be using `solve_ivp`, which is a part of `scipy.integrate`, to solve ODE systems related to epidemic models. You can read the documentation for `solve_ivp` at [https://docs.scipy.org/doc/scipy/reference/generated/scipy.integrate.solve\\_ivp.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.integrate.solve_ivp.html).

`solve_ivp` takes the ODE as a function, a tuple containing the start and end time, and an array with the initial conditions as arguments, and returns a bunch object containing the solution and other information. We can solve the following ODE system with the following code.

$$\begin{bmatrix} y_1(t) \\ y_2(t) \end{bmatrix}' = \begin{bmatrix} y_2(t) \\ \sin(t) - 5y_2(t) - y_1(t) \end{bmatrix} \quad (3.1)$$
$$y_1(0) = 0, \quad y_2(0) = 1, \quad t \in [0, 3\pi]$$

```
import numpy as np
from scipy.integrate import solve_ivp

# define the ode system as given in the problem
def ode(t,y):
    return np.array([y[1], np.sin(t) - 5*y[1] - y[0]])

# define the t0 and tf parameters
t0 = 0
tf = 3*np.pi

# define the initial conditions
y0 = np.array([0,1])

# solve the system
sol = solve_ivp(ode, (t0,tf), y0)

# Plot the system
```

```
import matplotlib.pyplot as plt

# plot y_1 against y_2
plt.plot(sol.y[0],sol.y[1])
plt.xlabel('$y_1$')
plt.ylabel('$y_2$')
plt.show()
```

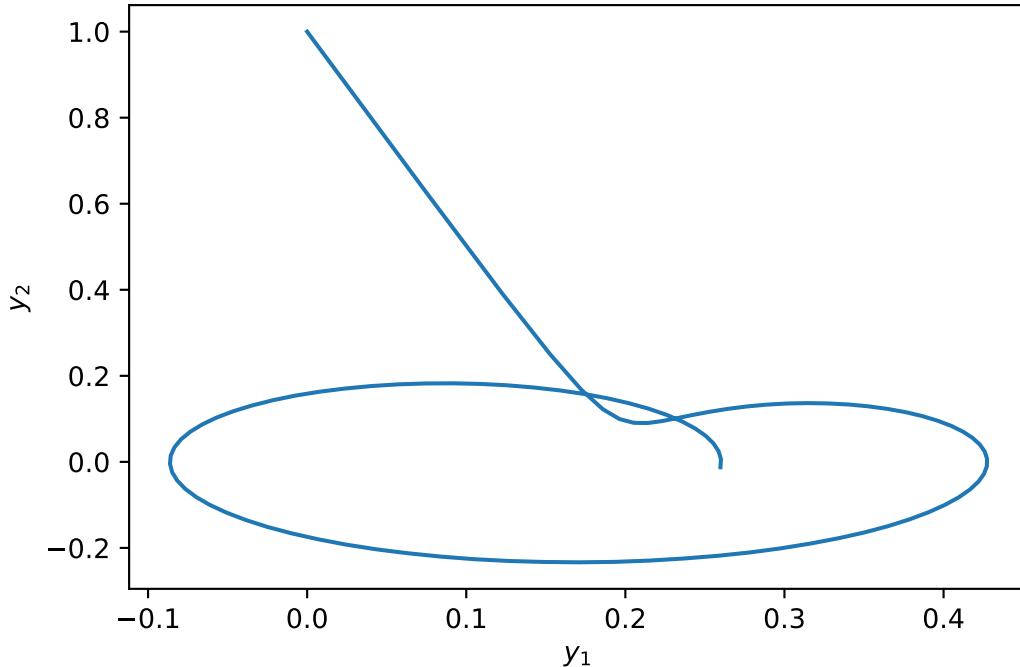


Figure 3.1: Solution to (3.1)

## The SIR Model

The SIR model describes the spread of an epidemic through a large population. It does this by describing the movement of the population through three phases of the disease: those individuals who are susceptible, those who are infectious, and those who have been removed from the disease. Those individuals in the removed class have either died, or have recovered from the disease and are now immune to it. If the outbreak occurs over a short period of time, we may reasonably assume that the total population is fixed, so that  $S'(t) + I'(t) + R'(t) = 0$ . We may also assume that  $S(t) + I(t) + R(t) = 1$ , so that  $S(t)$  represents the fraction of the population that is susceptible, etc.

Individuals may move from one class to another as described by the flow

$$S \rightarrow I \rightarrow R.$$

Let us consider the transition rate between  $S$  and  $I$ . Let  $\beta$  represent the average number of contacts made per unit time period (one day perhaps) that could spread the disease. The proportion of these contacts that are with a susceptible individual is  $S(t)$ . Thus, one infectious individual will on average infect  $\beta S(t)$  others per day. Let  $N$  represent the total population size. Then we obtain the differential equation

$$\frac{d}{dt}(S(t)N) = -\beta S(t)(I(t)N)$$

Now consider the transition rate between  $I$  and  $R$ . We assume that there is a fixed proportion  $\gamma$  of the infectious group who will recover on a given day, so that

$$\frac{d}{dt}R(t) = \gamma I(t).$$

Note that  $\gamma$  is the reciprocal of the average length of time spent in the infectious phase.

Since the derivatives sum to 0, we have  $I'(t) = -S'(t) - R'(t)$ , so the differential equations are given by

$$\frac{dS}{dt} = -\beta IS, \quad (3.2)$$

$$\frac{dI}{dt} = \beta IS - \gamma I, \quad (3.3)$$

$$\frac{dR}{dt} = \gamma I. \quad (3.4)$$

**Problem 1.** Suppose that, in a city of approximately three million, five people who have just become infectious have recently entered the city carrying a certain disease. Each of those individuals has one contact each day that could spread the disease, and an average of three days is spent in the infectious state. Find the solution of the corresponding SIR equations using `solve_ivp` for fifty days, where each time period is half a day, and plot your results. Use the percentages of each state, not the actual number of people in the state.

At the peak of the infection, how many in the city will still be able to work (assume for simplicity that those who are in the infectious state either cannot go to work or are unproductive, etc.)?

Hint: Use the `t`-values parameter in `solve_ivp` to pass in an array of `t`-values.

Compare your plot to Figure 1.

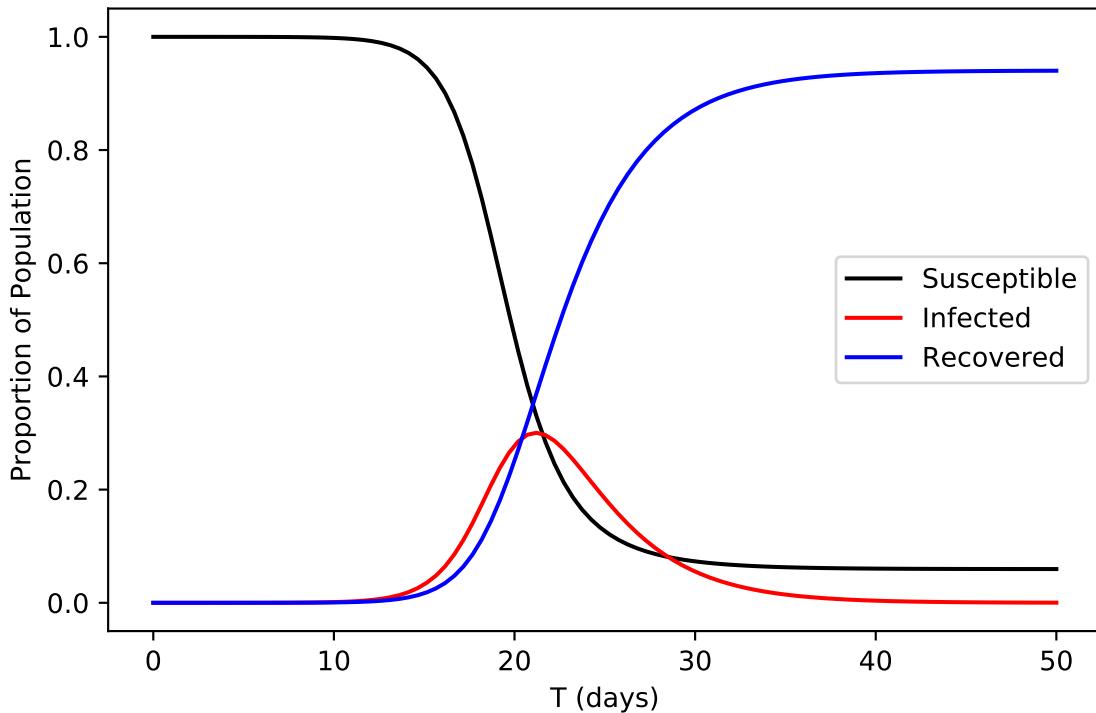


Figure 3.2: Solution to Problem (1)

SIR is an effective model for epidemic spread under certain assumptions. For example, we assume that the network is what's called "fully mixed." This implies that no group of members of a network are more likely to encounter each other than any other group. Because of this assumption, we should not use SIR to model networks we know to be poorly mixed. In fact, we should be clear in stating that almost no network is truly fully mixed; however this model is still effective for networks that are reasonably well mixed. In the next problem we will be using SIR to model data from the recent COVID-19 outbreak. To adhere to the "reasonably well mixed" criteria, we will be using only data from one county at a time.

**Problem 2.** On March 11, 2020, New York City had 52 confirmed cases of COVID-19. On that day New York started its lock-down measures. Using the following information, model what the spread of the virus could have been, using `solve_ivp()`, if New York did not implement any measures to curb the spread of the virus over the next 150 days:

- There are approximately 8.399 million people in New York city.
- The average case of COVID-19 lasts for 10 days.
- Each infected person can spread the virus to 2.5 people.

Plot your results for each day and compare to Figure 3.3.

1. At the projected peak, how many concurrent active cases are there?
2. Assuming that about 5% of COVID-19 cases require hospitalization, and using the fact that there are about 58,000 hospital beds in NYC, how many beds over capacity will the hospitals in NYC be at the projected peak?

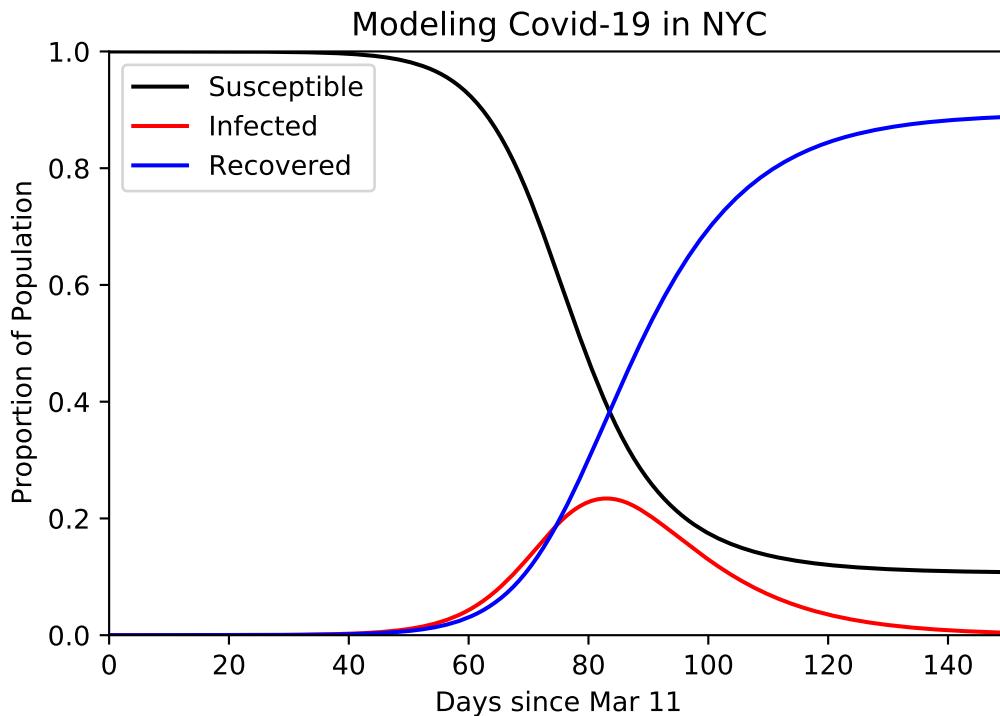


Figure 3.3: Solution to Problem (2).

## Variations on the SIR Model

The SIS model is a common variation of the SIR model. SIS Models describe diseases where individuals who have recovered from the disease do not gain any lasting immunity. There are only two compartments in this model: those who are susceptible, and those who are infectious. Here,  $f$  is the rate of becoming susceptible again.

The basic equations are given by

$$\begin{aligned}\frac{dS}{dt} &= -\beta IS + fI, \\ \frac{dI}{dt} &= \beta IS - fI\end{aligned}$$

Another alteration we can make to the SIR model is to add a birth and death rate. In the equations below we are assuming that the natural death rate together with the death rate caused by the disease is equal to the birth rate. This model is given by

$$\begin{aligned}\frac{dS}{dt} &= \mu(1 - S) - \beta IS, \\ \frac{dI}{dt} &= \beta IS - (\gamma + \mu)I, \\ \frac{dR}{dt} &= \gamma I - \mu R\end{aligned}$$

where  $\mu$  represents the death rate and equal birth rate, noting that any new person born is born into the susceptible population.

If we combine the last two variations we made on the SIR model we come to this formulation, which is an SIRS model. This SIRS model allows the transfer of individuals from the recovered/removed class to the susceptible class and includes modeling of the birth and death rates.

$$\frac{dS}{dt} = fR + \mu(1 - S) - \beta IS, \quad (3.5)$$

$$\frac{dI}{dt} = \beta IS - (\gamma + \mu)I, \quad (3.6)$$

$$\frac{dR}{dt} = -fR + \gamma I - \mu R. \quad (3.7)$$

**Problem 3.** There are 7 billion people in the world. Influenza, or the flu, is one of those viruses that everyone can be susceptible to, even after recovering. The flu virus is able to change in order to evade our immune system, and we become susceptible once more, although technically it is now a different strain. Suppose the virus originates with 1000 people in Texas after Hurricane Harvey flooded Houston, and stagnant water allowed the virus to proliferate. According to WebMD, once you get the virus, adults are contagious up to a week and kids up to 2 weeks. For this lab, suppose you are contagious for 10 days before recovering. Also suppose that on average someone makes one contact every two days that could spread the flu. Since we can catch a new strain of the flu, suppose that a recovered individual becomes susceptible again with probability  $f = 1/50$ . The flu is also known to be deadly, killing hundreds of thousands every year on top of the normal death rate. To assure a steady population, let the birth rate balance out the death rate, and in particular let  $\mu = .0001$ .

Using the SIRS model above, plot the proportion of population that is Susceptible, Infected, and Recovered over a one-year span (365 days). Compare your plot to Figure 3.4.

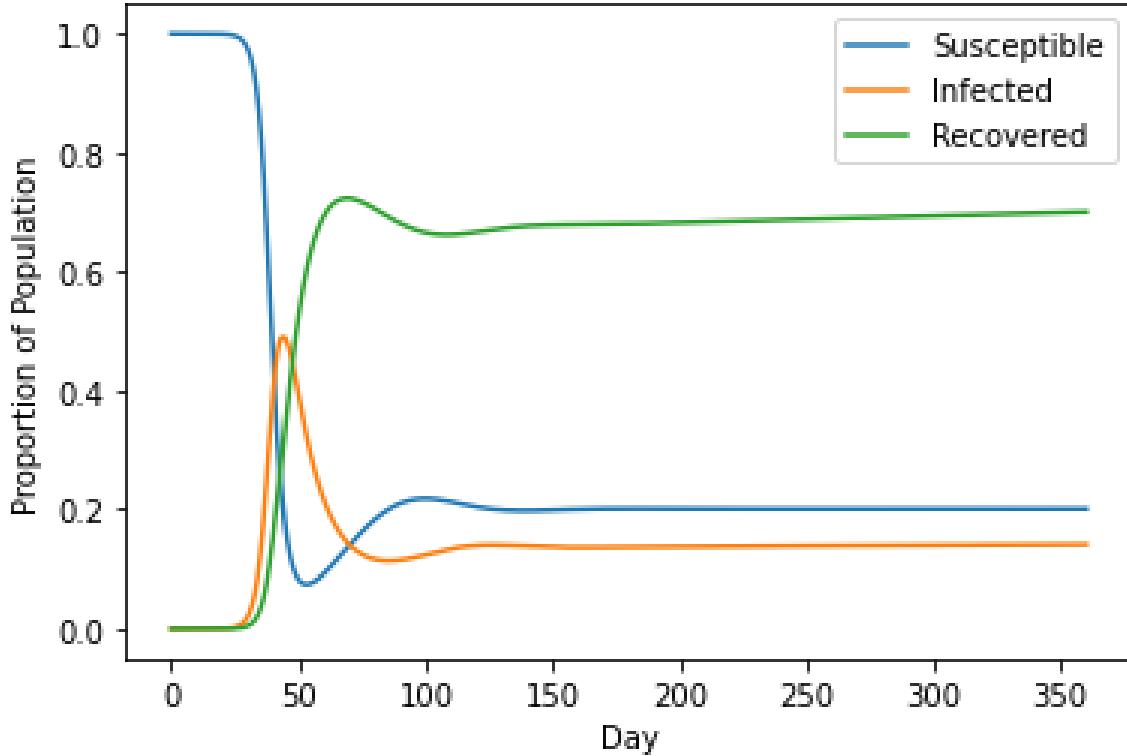


Figure 3.4: Solution to Problem (3).

### Modeling COVID-19 with Social Distancing

Social distancing upsets the main assumption that is made when trying to model epidemic spread using SIR models. During the periods of lockdown instituted by governments, the interaction networks between people in a city or county were disrupted to the point that standard SIR models were no longer effective at modeling the spread of COVID-19. A paper released in May of 2020 presented some alternative models for COVID-19 that have some success in modeling its spread during periods of social distancing.

This model claims that the growth of  $I(t)$  is polynomial with exponential decay (PGED). So we get the following form

$$I(t) \approx Bt^\alpha e^{-t/T_G},$$

which results in the following SIR type model

$$\frac{dS}{dt} = -\frac{\alpha}{t} I, \quad (3.8)$$

$$\frac{dI}{dt} = \left( \frac{\alpha}{t} - \frac{1}{T_G} \right) I, \quad (3.9)$$

$$\frac{dR}{dt} = \frac{1}{T_G} I, \quad (3.10)$$

where  $\alpha$  and  $T_G$  are simply model parameters. In this model  $\alpha T_G$  can be interpreted as the time of epidemic peak.

## Fitting Models

Model fitting can be a frustrating task if we only use our intuition and guess and check. Thankfully, SciPy's `optimize` library has tools we can use to make these problems a lot easier. Many of the functions in this library are designed to take an arbitrary function and find whatever input makes the output close to zero. Our job is to create a function that outputs zero at the right values.

Suppose we have some data that we believe to follow a cubic trend with the following model

$$\alpha x^3 + \beta(x^2 + 2x) + \delta.$$

In order to fit the data to this model we can use `scipy.optimize.minimize` and create a function that will output zero when the correct parameters are input. `scipy.optimize.minimize` will then return an `OptimizeResult` object, which contains the optimal parameters.

```
# import the minimizer function
from scipy.optimize import minimize

# load the data and get the x and y values
data = np.load('to_fit.npy')
xs = data[:,0]
ys = data[:,1]

# define the function we want to minimize
def fun(params):
    # unpack the parameters
    a,b,d = params

    # get the model output based on the parameters
    out = a*xs**3 + b*(xs**2 + 2*xs) + d

    # find the difference between out and the data
    diff = out - ys

    # must return a float
    return np.linalg.norm(diff)

# make a guess for the parameters
p0 = (1,1,1)

# find the best parameters for this model
minimize(fun,p0)
```

**Problem 4.** Fit the PEGD model to the COVID-19 data provided in `new_york_cases.npy`. Plot your results against  $1 - S(t)$ .

Hint: Set  $t_0 = 1$  as the PEGD model requires to divide by  $t$ , so we must have  $t \neq 0$ .

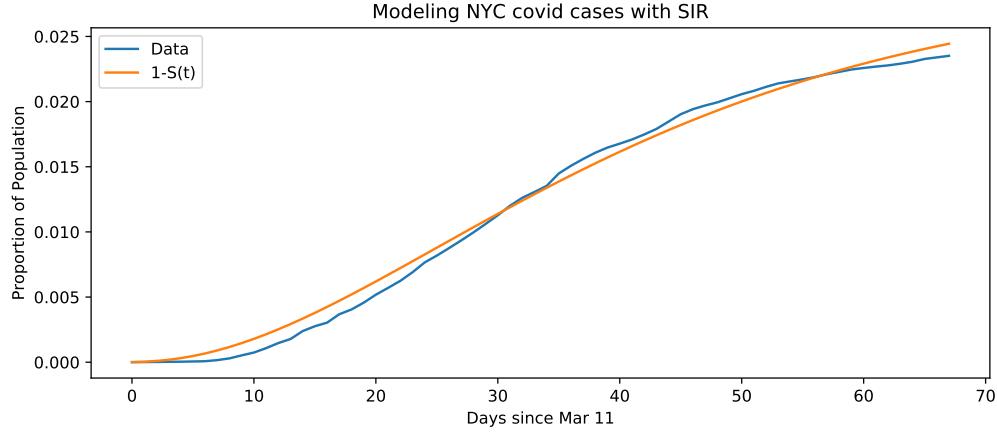


Figure 3.5: Solution to (4)

## Boundary Value Problems

The next exercise uses a variation of the SIR model called an SEIR model to describe the spread of measles<sup>1</sup>. This new model adds another compartment, called the exposed or latency phase. It assumes that the rate at which measles is contracted depends on the season, i.e. the rate is periodic. That allows us to formulate the yearly occurrence rate for measles as a boundary value problem. The boundary value problem looks like

$$\begin{bmatrix} S \\ E \\ I \end{bmatrix}' = \begin{bmatrix} \mu - \beta(t)SI \\ \beta(t)SI - E/\lambda \\ E/\lambda - I/\eta \end{bmatrix}, \quad (3.11)$$

$$\begin{aligned} S(0) &= S(1), \\ E(0) &= E(1), \\ I(0) &= I(1) \end{aligned} \quad (3.12)$$

Parameters  $\mu$  and  $\lambda$  represent the birth rate of the population and the latency period of measles, respectively.  $\eta$  represents the infectious period before an individual moves from the infectious class to the recovered class. After recovery an individual remains immune, which is why  $R(t)$  is not included in the system. The set up of this problem is not normal since we are excluding  $R(t)$ , but it results in a nice graph.

To solve this problem we will use a full-featured BVP solver that is available in SciPy. The code below demonstrates how to use `solve_bvp` to solve the BVP

$$\varepsilon y'' + yy' - y = 0, \quad y(-1) = 1, \quad y(1) = -1/3, \quad \varepsilon = .1 \quad (3.13)$$

Look at figure 3.6 for the solution.

---

<sup>1</sup>Numerical Solution of Boundary Value Problems for Ordinary Differential Equations, by Aescher, Mattheij, and Russell

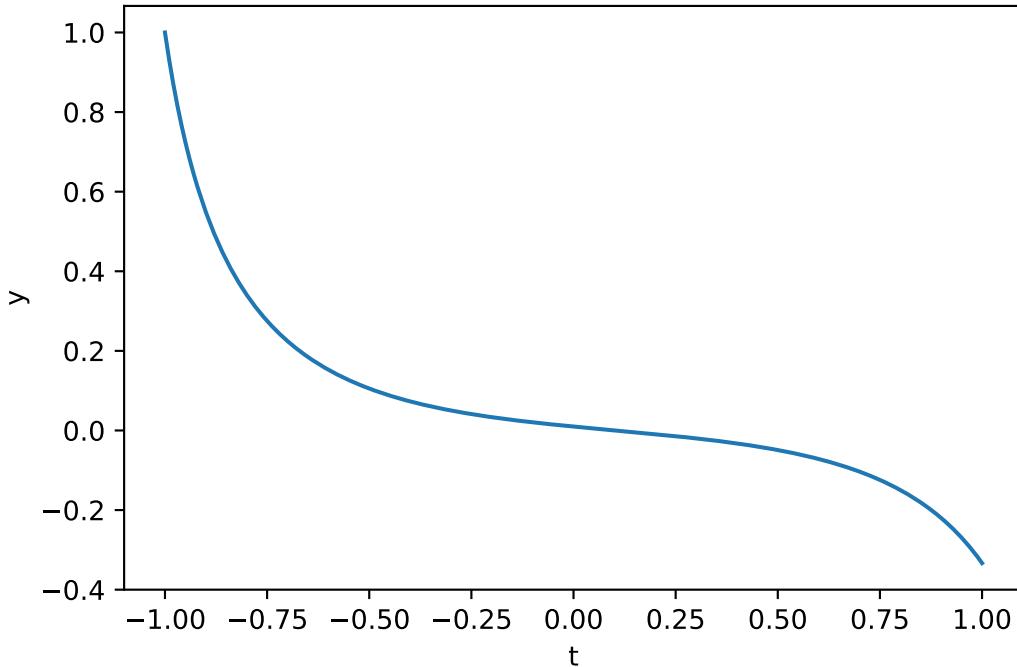


Figure 3.6: Solution to Equation (3.13)

The BVP solver expects you to pass it the boundary conditions as a callable function that computes the difference between a guess at the boundary conditions and the desired boundary conditions. When we use the BVP solver, we will tell it how many constraints there should be on each side of the domain so it knows how many entries to expect. In this case, we have one boundary condition on either side. These constraints are expected to evaluate to 0 when the boundary condition is satisfied.

```

import numpy as np
from scipy.integrate import solve_bvp
import matplotlib.pyplot as plt

epsilon, lbc, rbc = .1, 1, - 1/3

# The ode function takes the independent variable first
# It has return shape (n,)
def ode(x , y):
    return np.array([y[1] , (1/epsilon) * (y[0] - y[0] * y[1])])

# The return shape of bcs() is (n,)
def bcs(ya, yb):
    BCa = np.array([ya[0] - lbc])    # 1 Boundary condition on the left
    BCb = np.array([yb[0] - rbc])    # 1 Boundary condition on the right
    # The return values will be 0s when the boundary conditions are met exactly
    return np.hstack([BCa, BCb])

```

```

# The independent variable has size (m,) and goes from a to b with some step ←
# size
X = np.linspace(-1, 1, 200)
# The y input must have shape (n,m) and includes our initial guess for the ←
# boundaries
y = np.array([-1/3, -4/3]).reshape((-1,1))*np.ones((2, len(X)))

# There are multiple returns from solve_bvp(). We are interested in the y ←
# values which can be found in the sol field.
solution = solve_bvp(ode, bcs, X, y)
# We are interested in only y, not y', which is found in the first row of sol.
y_plot = solution.sol(X)[0]

plt.plot(X, y_plot)
plt.xlabel('t')
plt.ylabel('y')
plt.show()

```

**Problem 5.** Consider equations (3.11) and (3.12). Let the periodic function for our measles case be  $\beta(t) = \beta_0(1 + \beta_1 \cos 2\pi t)$ . Use parameters  $\beta_1 = 1$ ,  $\beta_0 = 1575$ ,  $\eta = 0.01$ ,  $\lambda = .0279$ , and  $\mu = .02$ . Note: in this case, time is measured in years, so run the solution over the interval  $[0, 1]$  to show a one-year cycle. The boundary conditions in (3.12) are just saying that the year will begin and end in the same state.

One issue that we encounter with this problem is that we have 6 boundary conditions but we only have 3 free variables. The 6 boundary conditions are the initial and final conditions of  $S$ ,  $E$ , and  $I$ . `solve_bvp` only allows as many boundary conditions as there are free variables, so what we can do is include “dummy” variables in the ODE. This allows more boundary conditions in the BVP solver, while not changing the ODE system that we are solving. To deal with this, let  $C(t) = [C_1(t), C_2(t), C_3(t)]$ , and add the equation

$$C'(t) = 0$$

to the system of ODEs given above (for a total of 6 equations) resulting in this final 6 variable system

$$\begin{bmatrix} S(t) \\ E(t) \\ I(t) \\ C_1 \\ C_2 \\ C_3 \end{bmatrix}' = \begin{bmatrix} \mu - \beta(t)SI \\ \beta(t)SI - E/\lambda \\ E/\lambda - I/\eta \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

We can then apply all 6 of the boundary conditions that we need. The boundary conditions can be separated using the following trick:

$$\begin{pmatrix} C_1(0) \\ C_2(0) \\ C_3(0) \end{pmatrix} = \begin{pmatrix} S(0) \\ E(0) \\ I(0) \end{pmatrix}, \quad \begin{pmatrix} C_1(1) \\ C_2(1) \\ C_3(1) \end{pmatrix} = \begin{pmatrix} S(1) \\ E(1) \\ I(1) \end{pmatrix}.$$

Now  $C_1, C_2, C_3$  become the 4th, 5th, and 6th rows of your solution matrix, so the 3 boundary conditions for the left are obtained by subtracting the last three entries of  $y(0)$  from the first three entries, giving you  $ya[0 : 3] - ya[3 :]$ . Similarly, your right boundary conditions will look like  $yb[0 : 3] - yb[3 :]$ .

When you code your boundary conditions, note that `solve_bvp` changes the initial conditions to force all the entries in the return of `bcs()` to be zero. You can use the initial conditions from Fig. 3.7 as your initial guess (which will be an array of 6 elements). Remember that the initial infected proportion is small, not 0.

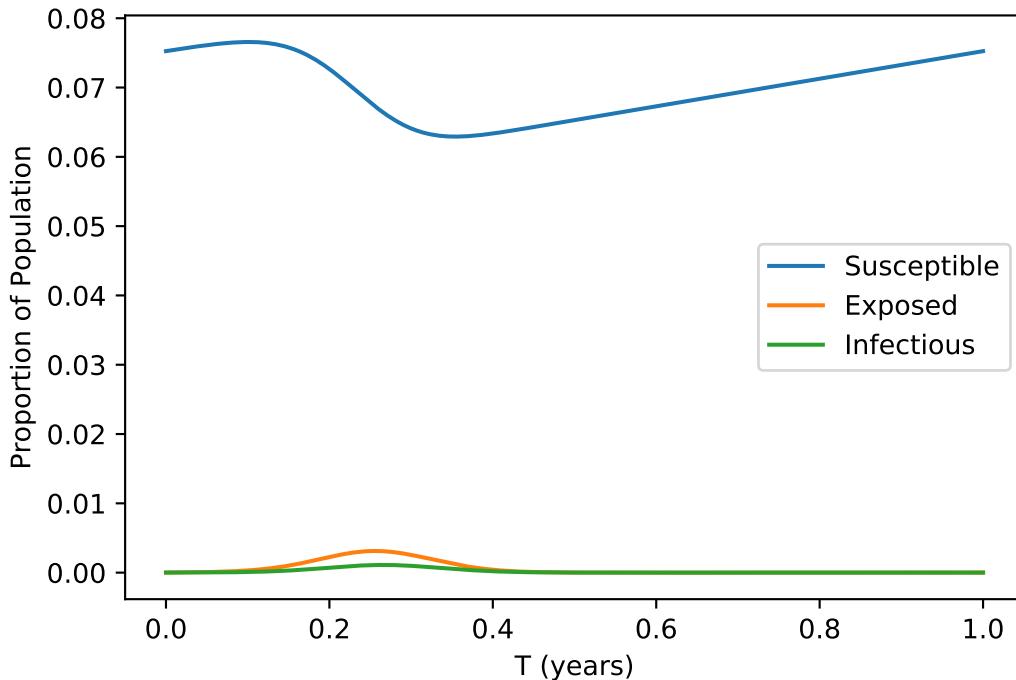


Figure 3.7: Solution to Problem (5)

# 4

# Numerical Methods for Initial Value Problems; Harmonic Oscillators

**Lab Objective:** Implement several basic numerical methods for initial value problems (IVPs) and use them to study harmonic oscillators.

## Methods for Initial Value Problems

Consider the initial value problem (IVP)

$$\begin{aligned}\mathbf{x}'(t) &= f(\mathbf{x}(t), t), \quad t_0 \leq t \leq t_f \\ \mathbf{x}(t_0) &= \mathbf{x}_0,\end{aligned}\tag{4.1}$$

where  $f$  is a suitably continuous function. A solution of (4.1) is a continuously differentiable, and possibly vector-valued, function  $\mathbf{x}(t) = [x_1(t), \dots, x_m(t)]^\top$ , whose derivative  $\mathbf{x}'(t)$  equals  $f(\mathbf{x}(t), t)$  for all  $t \in [t_0, t_f]$ , and for which the initial value  $\mathbf{x}(t_0)$  equals  $\mathbf{x}_0$ .

Under the right conditions, namely that  $f$  is uniformly Lipschitz continuous in  $\mathbf{x}(t)$  near  $\mathbf{x}_0$  and continuous in  $t$  near  $t_0$ , (4.1) is well-known to have a unique solution. However, for many IVPs, it is difficult, if not impossible, to find a closed-form, analytic expression for  $\mathbf{x}(t)$ . In these cases, numerical methods can be used to instead approximate  $\mathbf{x}(t)$ .

As an example, consider the initial value problem

$$\begin{aligned}x'(t) &= \sin(x(t)), \\ x(0) &= x_0.\end{aligned}\tag{4.2}$$

The solution  $x(t)$  is defined implicitly by

$$t = \ln \left| \frac{\cos(x_0) + \cot(x_0)}{\csc(x(t)) + \cot(x(t))} \right|.$$

This equation cannot be solved for  $x(t)$ , so it is difficult to understand what solutions to (4.2) look like. Since  $\sin(n\pi) = 0$ , there are constant solutions  $x_n(t) = n\pi$ ,  $n \in \mathbb{Z}$ . Using a numerical IVP solver, solutions for different values of  $x_0$  can be approximated. Figure 4.1 shows several of these approximate solutions, along with some of the constant, or equilibrium, solutions.

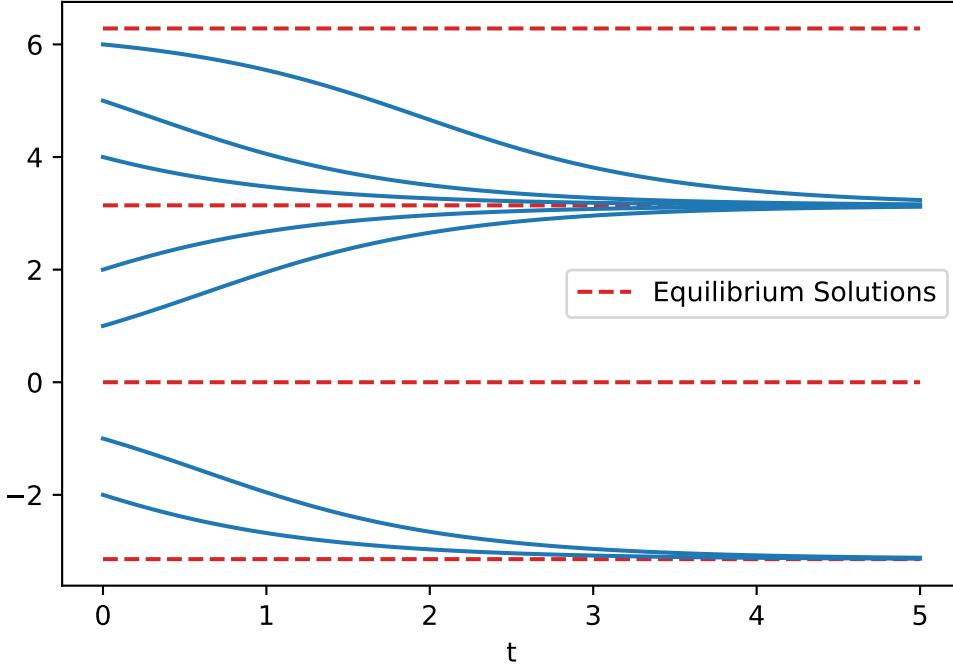


Figure 4.1: Several solutions of (4.2), using `scipy.integrate.odeint`.

## Numerical Methods

For the numerical methods that follow, the key idea is to seek an approximation for the values of  $\mathbf{x}(t)$  only on a finite set of values  $t_0 < t_1 < \dots < t_{n-1} < t_n (= t_f)$ . In other words, these methods try to solve for  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$  such that  $\mathbf{x}_i \approx \mathbf{x}(t_i)$ .

### Euler's Method

For simplicity, assume that each of the  $n$  subintervals  $[t_{i-1}, t_i]$  has equal length  $h = (t_f - t_0)/n$ .  $h$  is called the step size. Assuming  $\mathbf{x}(t)$  is twice-differentiable, for each component function  $x_j(t)$  of  $\mathbf{x}(t)$  and for each  $i$ , Taylor's Theorem says that

$$x_j(t_{i+1}) = x_j(t_i) + h x'_j(t_i) + \frac{h^2}{2} x''_j(c) \text{ for some } c \in [t_i, t_{i+1}].$$

The quantity  $\frac{h^2}{2} x''_j(c)$  is negligible when  $h$  is sufficiently small, and thus  $x_j(t_{i+1}) \approx x_j(t_i) + h x'_j(t_i)$ . Therefore, bringing the component functions of  $\mathbf{x}(t)$  back together gives

$$\begin{aligned} \mathbf{x}(t_{i+1}) &\approx \mathbf{x}(t_i) + h \mathbf{x}'(t_i), \\ &\approx \mathbf{x}(t_i) + h f(\mathbf{x}(t_i), t_i). \end{aligned}$$

This approximation leads to the Euler method: Starting with  $\mathbf{x}_0 = \mathbf{x}(t_0)$ ,  $\mathbf{x}_{i+1} = \mathbf{x}_i + h f(\mathbf{x}_i, t_i)$  for  $i = 0, 1, \dots, n - 1$ . Euler's method can be understood as starting with the point at  $\mathbf{x}_0$ , then calculating the derivative of  $\mathbf{x}(t)$  at  $t_0$  using  $f(\mathbf{x}_0, t_0)$ , followed by taking a step in the direction of the derivative scaled by  $h$ . Set that new point as  $\mathbf{x}_1$  and continue.

It is important to consider how the choice of step size  $h$  affects the accuracy of the approximation. Note that at each step of the algorithm, the local truncation error, which comes from neglecting the  $x_j''(c)$  term in the Taylor expansion, is proportional to  $h^2$ . The error  $\|\mathbf{x}(t_i) - \mathbf{x}_i\|$  at the  $i$ th step comes from  $i = \frac{t_i - t_0}{h}$  steps, which is proportional to  $h^{-1}$ , each contributing  $h^2$  error. Thus the global truncation error is proportional to  $h$ . Therefore, the Euler method is called a first-order method, or a  $\mathcal{O}(h)$  method. This means that as  $h$  gets small, the approximation of  $\mathbf{x}(t)$  improves in two ways. First,  $\mathbf{x}(t)$  is approximated at more values of  $t$  (more information about the solution), and second, the accuracy of the approximation at any  $t_i$  is improved proportional to  $h$  (better information about the solution).

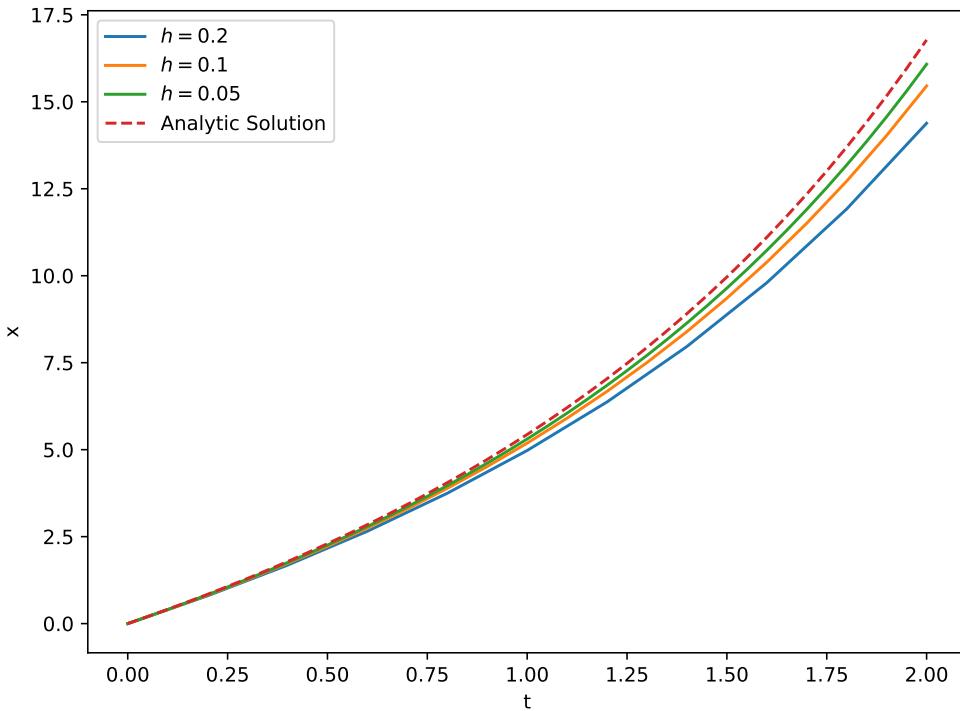


Figure 4.2: The solution of (4.3), alongside several approximations using Euler's method.

**Problem 1.** Write a function which implements Euler's method for an IVP of the form (4.1). Test your function on the IVP:

$$\begin{aligned} x'(t) &= x(t) - 2t + 4, \quad 0 \leq t \leq 2, \\ x(0) &= 0, \end{aligned} \tag{4.3}$$

where the analytic solution is  $x(t) = -2 + 2t + 2e^t$ . Use the Euler method to numerically approximate the solution with step sizes  $h = 0.2, 0.1$ , and  $0.05$ . Plot the true solution alongside the three approximations, and compare your results with Figure 4.2.

## Midpoint Method

The midpoint method is very similar to Euler's method. For small  $h$ , use the approximation

$$\mathbf{x}(t_{i+1}) \approx \mathbf{x}(t_i) + h f(\mathbf{x}(t_i) + \frac{h}{2} f(\mathbf{x}(t_i), t_i), t_i + \frac{h}{2}).$$

In this approximation, first set  $\hat{\mathbf{x}}_i = \mathbf{x}_i + \frac{h}{2} f(\mathbf{x}_i, t_i)$ , which is an Euler method step of size  $h/2$ . Then evaluate  $f(\hat{\mathbf{x}}_i, t_i + \frac{h}{2})$ , which is a more accurate approximation to the derivative  $\mathbf{x}'(t)$  in the interval  $[t_i, t_{i+1}]$ . Finally, a step is taken in that direction, scaled by  $h$ . It can be shown that the local truncation error for the midpoint method is  $\mathcal{O}(h^3)$ , giving global truncation error of  $\mathcal{O}(h^2)$ . This is a significant improvement over the Euler method. However, it comes at the cost of additional evaluations of  $f$  and a handful of extra floating point operations on the side. This tradeoff will be considered later in the lab.

## Runge-Kutta Methods

The Euler method and the midpoint method belong to a family called Runge-Kutta methods. There are many Runge-Kutta methods with varying orders of accuracy. Methods of order four or higher are most commonly used. A fourth-order Runge-Kutta method (RK4) iterates as follows:

$$\begin{aligned} K_1 &= f(\mathbf{x}_i, t_i), \\ K_2 &= f\left(\mathbf{x}_i + \frac{h}{2} K_1, t_i + \frac{h}{2}\right), \\ K_3 &= f\left(\mathbf{x}_i + \frac{h}{2} K_2, t_i + \frac{h}{2}\right), \\ K_4 &= f(\mathbf{x}_i + h K_3, t_{i+1}), \\ \mathbf{x}_{i+1} &= \mathbf{x}_i + \frac{h}{6}(K_1 + 2K_2 + 2K_3 + K_4). \end{aligned}$$

Runge-Kutta methods can be understood as a generalization of quadrature methods for approximating integrals, where the integrand is evaluated at specific points, and then the resulting values are combined in a weighted sum. For example, consider a differential equation

$$x'(t) = f(t)$$

Since the function  $f$  has no  $x$  dependence, this is a simple integration problem. In this case, Euler's method corresponds to the left-hand rule, the midpoint method becomes the midpoint rule, and RK4 reduces to Simpson's rule.

## Advantages of Higher-Order Methods

It can be useful to visualize the order of accuracy of a numerical method. A method of order  $p$  has relative error of the form

$$E(h) = Ch^p$$

taking the logarithm of both sides yields

$$\log(E(h)) = p \cdot \log(h) + \log(C)$$

Therefore, on a log-log plot against  $h$ ,  $E(h)$  is a line with slope  $p$  and intercept  $\log(C)$ .

**Problem 2.** Write functions that implement the midpoint and fourth-order Runge-Kutta methods. Use the Euler, Midpoint, and RK4 methods to approximate the value of the solution for the IVP (4.3) from Problem 1 for step sizes of  $h = 0.2, 0.1, 0.05, 0.025$ , and  $0.0125$ .

Plot the following graphs

- The true solution alongside the approximation obtained from each method when  $h = 0.2$ .
- A log-log plot (use `plt.loglog`) of the relative error  $|x(2) - x_n|/|x(2)|$  as a function of  $h$  for each approximation.

Compare your second plot with Figure 4.3.

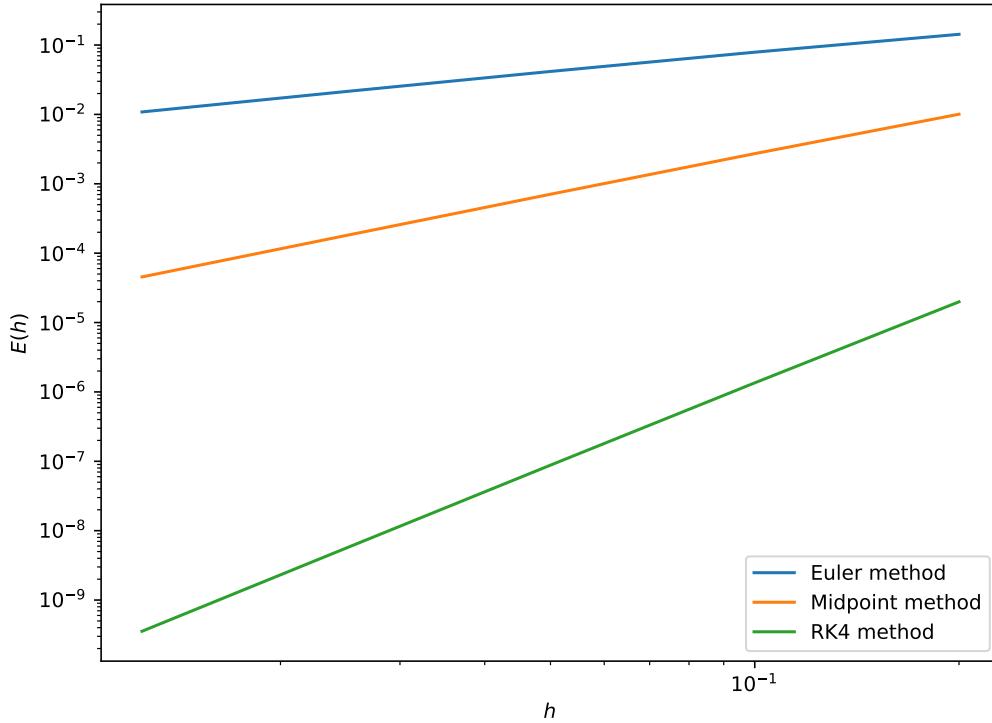


Figure 4.3: Loglog plot of the relative error in approximating  $x(2)$ , using step sizes  $h = 0.2, 0.1, 0.05, 0.025$ , and  $0.0125$ . The slope of each line demonstrates the first, second, and fourth order convergence of the Euler, Midpoint, and RK4 methods, respectively.

The Euler, midpoint, and RK4 methods help illustrate the potential trade-off between order of accuracy and computational expense. To increase the order of accuracy, more evaluations of  $f$  must be performed at each step. It is possible that this trade-off could make higher-order methods undesirable, as (in theory) one could use a lower-order method with a smaller step size  $h$ . However, this is not generally the case. Assuming efficiency is measured in terms of the number of  $f$ -evaluations required to reach a certain threshold of accuracy, higher-order methods turn out to be much more efficient. For example, consider the IVP

$$\begin{aligned} x'(t) &= x(t) \cos(t), \quad t \in [0, 8], \\ x(0) &= 1. \end{aligned} \tag{4.4}$$

Figure 4.4 illustrates the comparative efficiency of the Euler, Midpoint, and RK4 methods applied to (4.4). The higher-order RK4 method requires fewer  $f$ -evaluations to reach the same level of relative error as the lower-order methods. As  $h$  becomes small, which corresponds to increasing functional evaluations, each method reaches a point where the relative error  $|x(8) - x_n|/|x(8)|$  stops improving. This occurs when  $h$  is so small that floating point round-off error overwhelms local truncation error. Notice that the higher-order methods are able to reach a better level of relative error before this phenomena occurs.

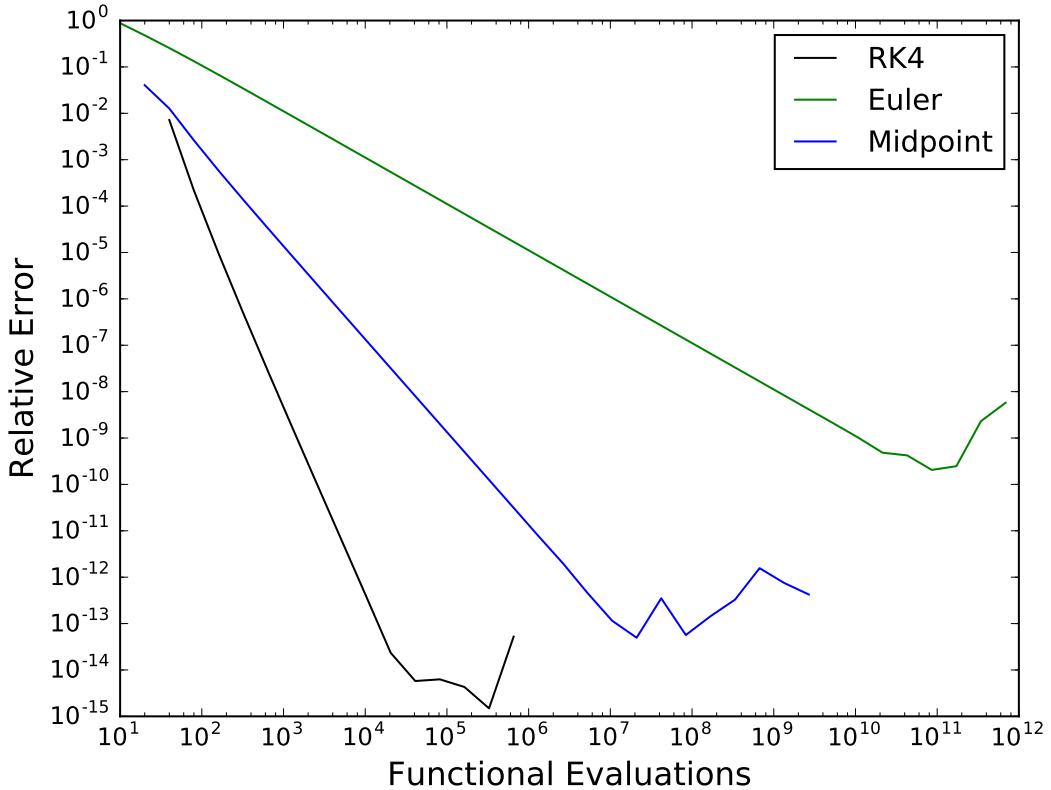


Figure 4.4: The relative error in computing the solution of (4.4) at  $x = 8$  versus the number of times the right-hand side of (4.4) must be evaluated.

## Harmonic Oscillators and Resonance

Harmonic oscillators are common in classical mechanics. A few examples include the pendulum (with small displacement), spring-mass systems, and the flow of electric current through various types of circuits. A harmonic oscillator  $y(t)$ <sup>1</sup> is a solution to an initial value problem of the form

$$\begin{aligned} my'' + \gamma y' + ky &= f(t), \\ y(0) = y_0, \quad y'(0) &= y'_0. \end{aligned}$$

Here,  $m$  represents the mass on the end of a spring,  $\gamma$  represents the effect of damping on the motion,  $k$  is the spring constant, and  $f(t)$  is the external force applied.

### Simple harmonic oscillators

A simple harmonic oscillator is a harmonic oscillator that is not damped,  $\gamma = 0$ , and is free,  $f = 0$ , rather than forced,  $f \neq 0$ . A simple harmonic oscillator can be described by the IVP

$$\begin{aligned} my'' + ky &= 0, \\ y(0) = y_0, \quad y'(0) &= y'_0. \end{aligned}$$

The solution of this IVP is  $y = c_1 \cos(\omega_0 t) + c_2 \sin(\omega_0 t)$ , where  $\omega_0 = \sqrt{k/m}$  is the natural frequency of the oscillator and  $c_1$  and  $c_2$  are determined by applying the initial conditions.

To solve this IVP using a Runge-Kutta method, it must be written in the form

$$\mathbf{x}'(t) = f(\mathbf{x}(t), t)$$

This can be done by setting  $x_1 = y$  and  $x_2 = y'$ . Then we have

$$\mathbf{x}' = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}' = \begin{bmatrix} x_2 \\ \frac{-k}{m}x_1 \end{bmatrix}$$

Therefore

$$f(\mathbf{x}(t), t) = \begin{bmatrix} x_2 \\ \frac{-k}{m}x_1 \end{bmatrix}$$

**Problem 3.** Use the RK4 method to solve the simple harmonic oscillator satisfying

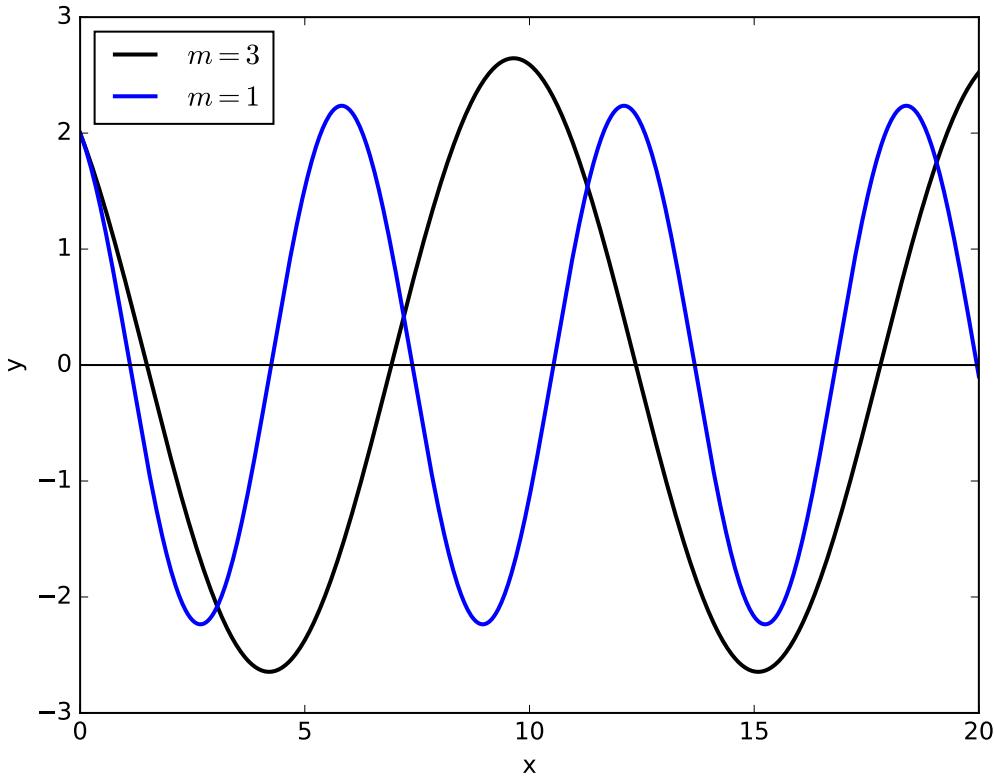
$$\begin{aligned} my'' + ky &= 0, \quad 0 \leq t \leq 20, \\ y(0) = 2, \quad y'(0) &= -1, \end{aligned} \tag{4.5}$$

for  $m = 1$  and  $k = 1$ .

Plot your numerical approximation of  $y(t)$ . Compare this with the numerical approximation when  $m = 3$  and  $k = 1$ . Consider: Why does the difference in solutions make sense physically?

---

<sup>1</sup>It is customary to write  $y$  instead of  $y(t)$  when it is unambiguous that  $y$  denotes the dependent variable.

Figure 4.5: Solutions of (4.5) for several values of  $m$ .

## Damped free harmonic oscillators

A damped free harmonic oscillator  $y(t)$  satisfies the IVP

$$\begin{aligned} my'' + \gamma y' + ky &= 0, \\ y(0) = y_0, \quad y'(0) &= y'_0. \end{aligned}$$

The roots of the characteristic equation are

$$r_1, r_2 = \frac{-\gamma \pm \sqrt{\gamma^2 - 4km}}{2m}.$$

Note that the real parts of  $r_1$  and  $r_2$  are always negative, and so any solution  $y(t)$  will decay over time due to a dissipation of the system energy. There are several cases to consider for the general solution of this equation:

1. If  $\gamma^2 > 4km$ , then the general solution is  $y(t) = c_1 e^{r_1 t} + c_2 e^{r_2 t}$ . Here the system is said to be overdamped. Notice from the general solution that there is no oscillation in this case.
2. If  $\gamma^2 = 4km$ , then the general solution is  $y(t) = c_1 e^{\gamma t/2m} + c_2 t e^{\gamma t/2m}$ . Here the system is said to be critically damped.

3. If  $\gamma^2 < 4km$ , then the general solution is

$$\begin{aligned} y(t) &= e^{-\gamma t/2m} [c_1 \cos(\mu t) + c_2 \sin(\mu t)], \\ &= Re^{-\gamma t/2m} \sin(\mu t + \delta), \end{aligned}$$

where  $R$  and  $\delta$  are fixed, and  $\mu = \sqrt{4km - \gamma^2}/2m$ . This system does oscillate.

**Problem 4.** Use the RK4 method to solve for the damped free harmonic oscillator satisfying

$$\begin{aligned} y'' + \gamma y' + y &= 0, \quad 0 \leq t \leq 20, \\ y(0) &= 1, \quad y'(0) = -1. \end{aligned}$$

For  $\gamma = 1/2$ , and  $\gamma = 1$ , simultaneously plot your numerical approximations of  $y$ .

## Forced harmonic oscillators without damping

Consider the systems described by the differential equation

$$my''(t) + ky(t) = F(t). \quad (4.6)$$

In many instances, the external force  $F(t)$  is periodic, so assume that  $F(t) = F_0 \cos(\omega t)$ . If  $\omega_0 = \sqrt{k/m} \neq \omega$ , then the general solution of 4.6 is given by

$$y(t) = c_1 \cos(\omega_0 t) + c_2 \sin(\omega_0 t) + \frac{F_0}{m(\omega_0^2 - \omega^2)} \cos(\omega t).$$

If  $\omega_0 = \omega$ , then the general solution is

$$y(t) = c_1 \cos(\omega_0 t) + c_2 \sin(\omega_0 t) + \frac{F_0}{2m\omega_0} t \sin(\omega_0 t).$$

When  $\omega_0 = \omega$ , the solution contains a term that grows arbitrarily large as  $t \rightarrow \infty$ . If we included damping, then the solution would be bounded but large for small  $\gamma$  and  $\omega$  close to  $\omega_0$ .

Consider a physical spring-mass system. Equation 4.6 holds only for small oscillations; this is where Hooke's law is applicable. However, the fact that the equation predicts large oscillations suggests the spring-mass system could fall apart as a result of the external force. This mechanical resonance has been known to cause failure of bridges, buildings, and airplanes.

**Problem 5.** Use the RK4 method to solve the damped and forced harmonic oscillator satisfying

$$\begin{aligned} 2y'' + \gamma y' + 2y &= 2 \cos(\omega t), \quad 0 \leq t \leq 40, \\ y(0) &= 2, \quad y'(0) = -1. \end{aligned} \quad (4.7)$$

For the following values of  $\gamma$  and  $\omega$ , plot your numerical approximations of  $y(t)$ :  $(\gamma, \omega) = (0.5, 1.5)$ ,  $(0.1, 1.1)$ , and  $(0, 1)$ . Compare your results with Figure 4.7.

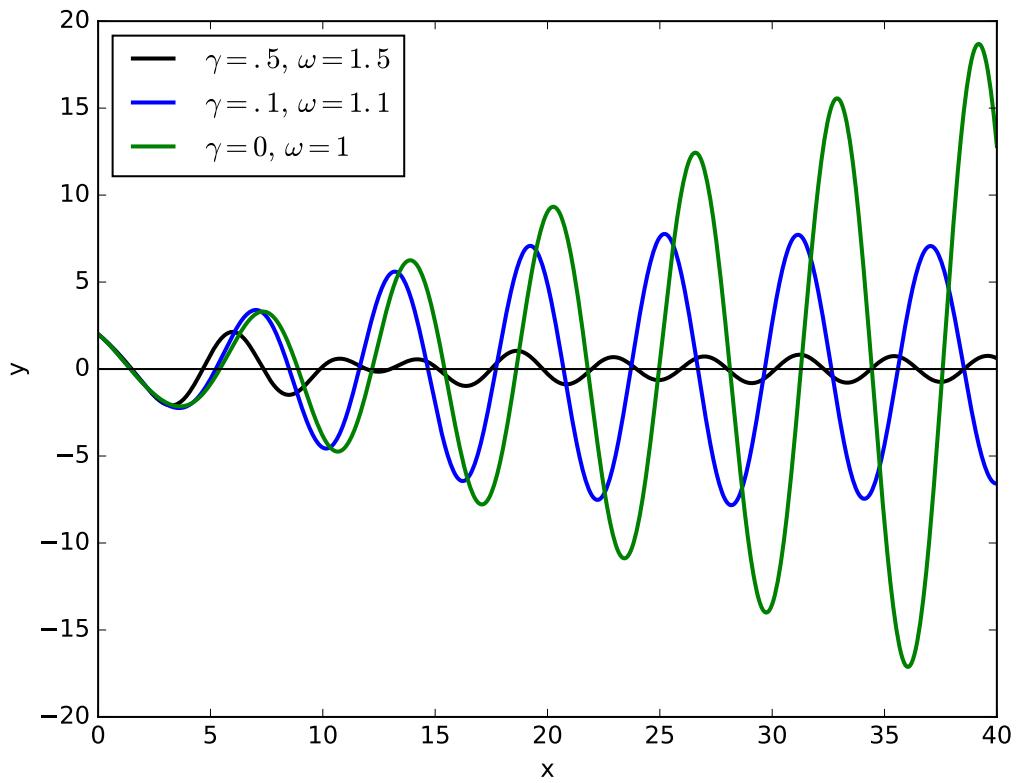


Figure 4.6: Solutions of (4.7) for several values of  $\omega$  and  $\gamma$ .

# 5

# Predator-Prey and Weight Change Models

**Lab Objective:** We introduce built-in methods for solving Initial Value Problems and apply the methods to two dynamical systems. The first system looks at the relationship between a predator and its prey. The second model is a weight change model based on thermodynamics and kinematics.

## Predator-Prey Model

ODEs are commonly used to model relationships between predator and prey populations. For example, consider the populations of wolves, the predator, and rabbits, the prey, in Yellowstone National Park. Let  $r(t)$  and  $w(t)$  represent the rabbit and wolf populations respectively at time  $t$ , measured in years. We will make a few assumptions to simplify our model:

- In the absence of wolves, the rabbit population grows at a positive rate proportional to the current population. Thus when  $w(t) = 0$ ,  $dr/dt = \alpha r(t)$ , where  $\alpha > 0$ .
- In the absence of rabbits, the wolves die out. Thus when  $r(t) = 0$ ,  $dw/dt = -\delta w(t)$ , where  $\delta > 0$ .
- The number of encounters between rabbits and wolves is proportional to the product of their populations. The wolf population grows proportional to the number of encounters by  $\beta r(t)w(t)$  (where  $\beta > 0$ ), and the rabbit population decreases proportional to the number of encounters by  $-\gamma r(t)w(t)$  (where  $\gamma > 0$ ).

This leads to the following system of ODEs:

$$\begin{aligned}\frac{dr}{dt} &= \alpha r - \beta rw = r(\alpha - \beta w) \\ \frac{dw}{dt} &= -\delta w + \gamma rw = w(-\delta + \gamma r)\end{aligned}\tag{5.1}$$

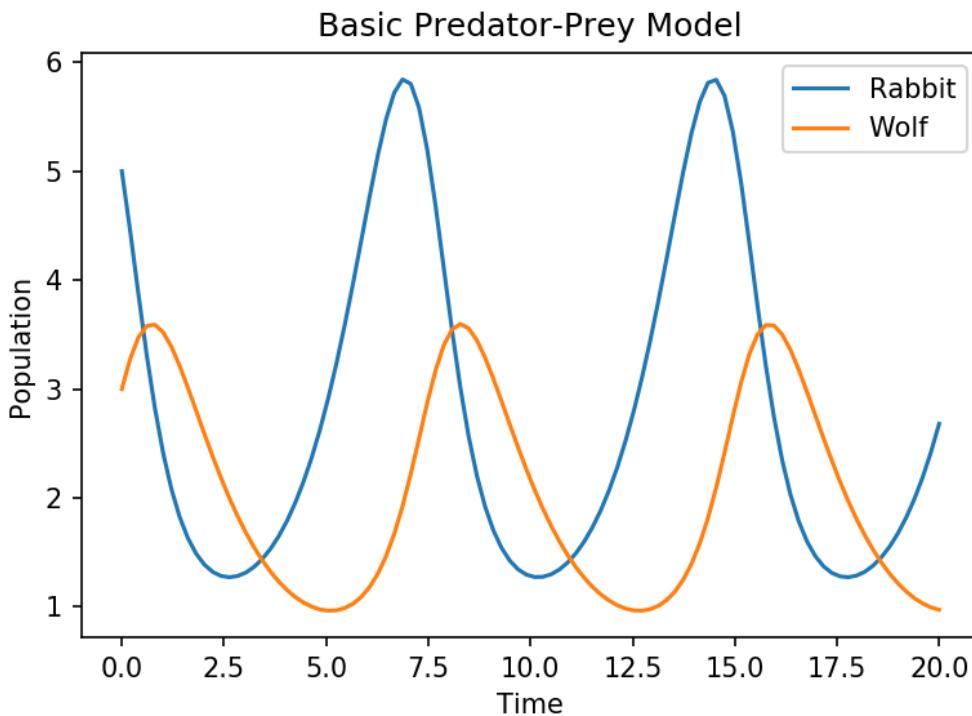


Figure 5.1: The solution to the system found in (5.1)

**Problem 1.** Define the function `predator_prey()` that accepts the current  $r(t)$  and  $w(t)$  values as a 1d array  $y$ , and the current time  $t$ , and returns the right hand side of (5.1) as an ndarray. Use  $\alpha = 1.0$ ,  $\beta = 0.5$ ,  $\delta = 0.75$ , and  $\gamma = 0.25$  as your growth parameters.

Hint: you will want to use `solve_ivp`.

**Problem 2.** Use `solve_ivp` to solve (5.1) with initial conditions  $(r_0, w_0) = (5, 3)$  and time ranging from 0 to 20 years. Display the resulting rabbit and wolf populations over time (stored as rows in the attribute `y` of the output of `solve_ivp`) on the same plot. Your graph should match the graph in figure 5.1.

## Variations on the Predator-Prey

### The Lotka-Volterra model

The representation of the predator-prey relationship found in (5.1) is called the Lotka-Volterra predator-prey model and is typically given by

$$\begin{aligned}\frac{du}{dt} &= \alpha u - \beta uv, \\ \frac{dv}{dt} &= -\delta v + \gamma uv.\end{aligned}$$

where  $u$  and  $v$  represent the prey and predator populations, respectively. Here  $\alpha$ ,  $\beta$ ,  $\delta$ , and  $\gamma$  are the same as before but now for an arbitrary prey and predator.

The equilibria (fixed points) of a system occur when the derivatives are zero. In this example, that occurs at  $(u, v) = (0, 0)$  and  $(u, v) = (\frac{c}{d}, \frac{a}{b})$ . Visualizing the phase portrait helps to give more insight into the dynamics of a system. We will do this by first nondimensionalizing our system to reduce the number of parameters. Let  $U = \frac{\gamma}{\delta}u$ ,  $V = \frac{\beta}{\alpha}v$ ,  $\bar{t} = \alpha t$ , and  $\eta = \frac{\gamma}{\alpha}$ . Substituting into the original ODEs we obtain the nondimensional system of equations

$$\begin{aligned}\frac{dU}{d\bar{t}} &= U(1 - V), \\ \frac{dV}{d\bar{t}} &= \eta V(U - 1).\end{aligned}\tag{5.2}$$

**Problem 3.** Similar to problem 1, define the function `Lotka_Volterra()` that takes in the current predator and prey populations as a 1d array  $y$  and the current time as a float  $t$  and returns the right hand side of the system (5.2) with  $\eta = 1/3$ .

The following three lines of code plot the phase portrait of (5.2). For more documentation on quiver plots see the documentation.

```
Y1, Y2 = np.meshgrid(np.linspace(0, 4.5, 25), np.linspace(0, 4.5, 25))
dU, dV = Lotka_Volterra(0, (Y1, Y2))
Q = plt.quiver(Y1[::3, ::3], Y2[::3, ::3], U[::3, ::3], V[::3, ::3])
```

Using `solve_ivp`, solve (5.2) with three different initial conditions  $y_0 = (1/2, 1/3)$ ,  $y_0 = (1/2, 3/4)$ , and  $y_0 = (1/16, 3/4)$  and time domain  $t = [0, 13]$ . Plot these three solutions on the same graph as the phase portrait and the equilibria  $(0, 0)$  and  $(1, 1)$ .

Since your solutions are being plotted with the phase portrait, plot the two populations against each other (instead of both individually against time). Your plot should match 5.2.

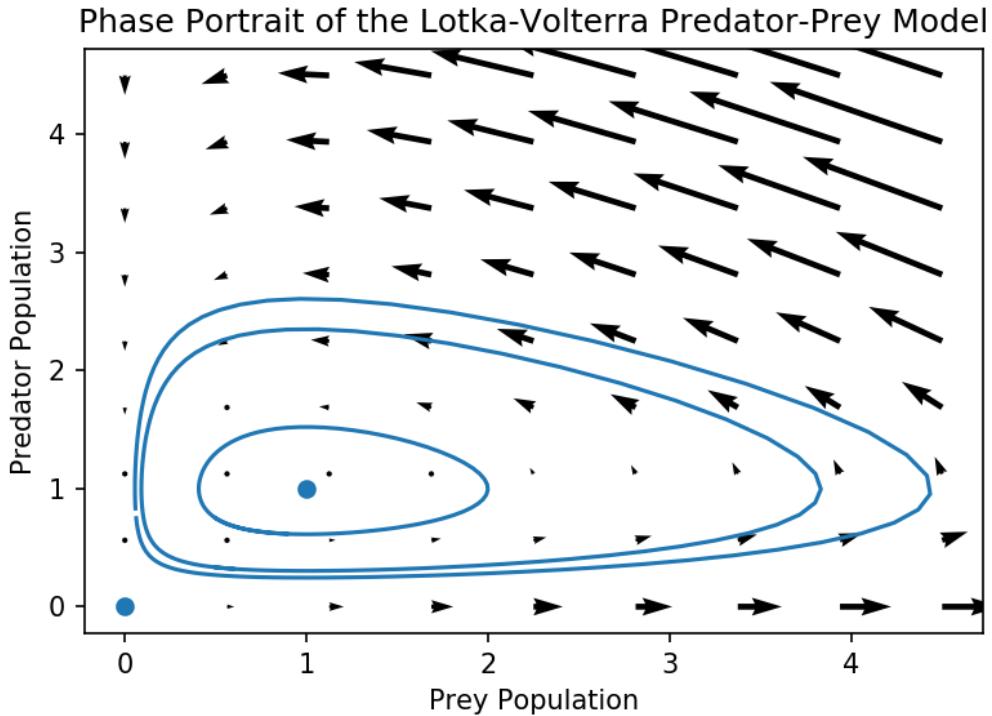


Figure 5.2: The phase portrait for the nondimensionalized Lotka-Volterra predator-prey equations with parameters  $\eta = 1/3$ .

### The Logistic model

Notice that the Lotka-Volterra equations predict prey populations will grow exponentially in the absence of predators. The logistic predator-prey equations change this dynamic by adding a carrying capacity  $K$  to the prey population:

$$\begin{aligned}\frac{du}{dt} &= \alpha u \left(1 - \frac{u}{K}\right) - \beta uv, \\ \frac{dv}{dt} &= -\delta v + \gamma uv.\end{aligned}$$

We can again do dimensional analysis on this system to simplify parameters. Let  $U = \frac{u}{K}$ ,  $V = \frac{\beta}{\alpha}v$ ,  $\bar{t} = \alpha t$ ,  $\eta = \frac{\gamma K}{\alpha}$ , and  $\rho = \frac{\delta}{\gamma K}$ . Then the nondimensional logistic equations are

$$\begin{aligned}\frac{dU}{d\bar{t}} &= U(1 - U - V), \\ \frac{dV}{d\bar{t}} &= \eta V(U - \rho).\end{aligned}\tag{5.3}$$

**Problem 4.** Define a new function `Logistic_Model()` that takes in the current predator and prey populations  $y$  and the current time  $t$  and returns the right hand side of (5.3) as a tuple. Use `solve_ivp` to compute solutions  $(U, V)$  of (5.3) for initial conditions  $(1/3, 1/3)$  and  $(1/2, 1/5)$  with  $(t_0, t_f) = (0, 13)$ . Do this for parameter values  $\eta, \rho = 1, 0.3$  and also for values  $\eta, \rho = 1, 1.1$ .

Create a phase portrait for the logistic equations using both sets of parameter values. Plot the direction field, all equilibrium points, and both solution orbits on the same plot for each set of parameter values.

## A Weight Change Model

The main idea behind weight change is simple. If a person takes in more energy than they expend, they gain weight. If they take in less than they expend, they lose weight. Let energy balance  $EB$  be the difference between energy intake  $EI$  and energy expenditure  $EE$ , so that

$$EB = EI - EE.$$

If the balance is positive, weight is gained and similarly if the balance is negative, weight is lost.

A person's body weight at a time  $t$  can be expressed as the sum of the weight of their fat tissue  $F(t)$  and the weight of their lean tissue  $L(t)$ ; that is,  $BW(t) = F(t) + L(t)$ . Using this, the change in body weight can be expressed as the following system of ODEs:

$$\begin{aligned} \frac{dF}{dt} &= \frac{(1 - p(t))EB(t)}{\rho_F}, \\ \frac{dL}{dt} &= \frac{p(t)EB(t)}{\rho_L}, \end{aligned} \tag{5.4}$$

where  $(1 - p(t))$  and  $p(t)$  represent the proportion of the energy balance ( $EB(t)$ ) that results in a change in the quantity of fatty or lean tissue, respectively. The constants  $\rho_F$  and  $\rho_L$  represent the energy density of fatty and lean tissue, approximated as  $\rho_F = 9400$  kcal/kg and  $\rho_L = 1800$  kcal/kg.

To solve this system, we first need to express  $p(t)$  and  $EB(t)$  in terms of  $F$  and  $L$ . These functions will also depend on physical activity level,  $PAL$ , and energy intake,  $EI$ , which vary among individuals.

We will find an expression for  $p(t)$  using Forbes' Law<sup>1</sup> which states that

$$\frac{dF}{dL} = \frac{F}{10.4}.$$

Notice

$$\frac{F}{10.4} = \frac{dF}{dL} = \frac{dF/dt}{dL/dt} = \frac{\frac{(1 - p(t))EB(t)}{\rho_F}}{\frac{p(t)EB(t)}{\rho_L}} = \frac{\rho_L}{\rho_F} \frac{1 - p(t)}{p(t)}.$$

Solving for  $p(t)$  gives Forbes' equation

$$p(t) = \frac{C}{C + F(t)} \quad \text{where} \quad C = 10.4 \frac{\rho_L}{\rho_F}. \tag{5.5}$$

---

<sup>1</sup>Lean body mass-body fat interrelationships in humans, Forbes, G.B.; Nutrition reviews, pgs 225-231, 1987.

We will now find an expression for  $EB(t)$ . Recall  $EB(t) = EI - EE$ . We will use the following expression for energy expenditure ( $EE$ ) to define  $EB(t)$ .

$$EE = PAL \times RMR \quad (5.6)$$

where  $PAL$  is physical activity level (as previously mentioned) and  $RMR$  is resting metabolic rate. Physical activity level can be determined using table 1.1.

1.40–1.69	People who are sedentary and do not exercise regularly, spend most of their time sitting, standing, with little body displacement
1.70–1.99	People who are active, with frequent body displacement throughout the day or who exercise frequently
2.00–2.40	People who engage regularly in strenuous work or exercise for several hours each day

Table 5.1: This is a rough guide for physical activity level (PAL).

We will use the following equation for computing  $RMR$ ,

$$RMR = K + \gamma_F F(t) + \gamma_L L(t) + \eta_F \frac{dF}{dt} + \eta_L \frac{dL}{dt} + \beta_{at} EI, \quad (5.7)$$

where  $\gamma_F = 3.2 \text{ kcal/kg/d}$ ,  $\gamma_L = 22 \text{ kcal/kg/d}$ ,  $\eta_F = 180 \text{ kcal/kg}$ , and  $\eta_L = 230 \text{ kcal/kg}^2$ <sup>3</sup>. Further, we let  $\beta_{at} = 0.14$  denote the coefficient for adaptive thermogenesis. Finally, we remark that the constant  $K$  can be tuned to an individual's body type directly through RMR and fat measurement, and is assumed to remain constant over time.

Thus, since the input  $EI$  is assumed to be known, we can use (5.6), (5.7) and (5.5) to write (5.4) in terms of  $F$  and  $L$ , thus allowing us to close the system of ODEs.

Specifically, we have

$$\begin{aligned} RMR &= \frac{EE}{PAL} = K + \gamma_F F(t) + \gamma_L L(t) + \eta_F \frac{dF}{dt} + \eta_L \frac{dL}{dt} + \beta_{at} EI \\ \frac{1}{PAL} (EE - EI + EI) &= K + \gamma_F F(t) + \gamma_L L(t) \\ &\quad + \left( \frac{\eta_F}{\rho_F} (1 - p(t)) + \frac{\eta_L}{\rho_L} p(t) \right) EB(t) + \beta_{at} EI. \\ \left( \frac{1}{PAL} - \beta_{at} \right) EI &= K + \gamma_F F(t) + \gamma_L L(t) \\ &\quad + \left( \frac{\eta_F}{\rho_F} (1 - p(t)) + \frac{\eta_L}{\rho_L} p(t) + \frac{1}{PAL} \right) EB(t). \end{aligned}$$

Solving for  $EB(t)$  in the last equation yields

$$EB(t) = \frac{\left( \frac{1}{PAL} - \beta_{at} \right) EI - K - \gamma_F F(t) - \gamma_L L(t)}{\frac{\eta_F}{\rho_F} (1 - p(t)) + \frac{\eta_L}{\rho_L} p(t) + \frac{1}{PAL}}. \quad (5.8)$$

<sup>2</sup>Modeling weight-loss maintenance to help prevent body weight regain; Hall, K.D. and Jordan, P.N.; The American journal of clinical nutrition, pg 1495, 2008

<sup>3</sup>Quantification of the effect of energy imbalance on bodyweight; Hall, K.D. et al.; The Lancet, pgs 826-837, 2011

In equilibrium ( $EB = 0$ ), this gives us

$$K = \left( \frac{1}{PAL} - \beta_{at} \right) EI - \gamma_F F - \gamma_L L. \quad (5.9)$$

Thus, for a subject who has maintained the same weight for a while, one can determine  $K$  by using (5.9), if they know their average caloric intake and amount of fat (assume  $L = BW - F$ ).

**Problem 5.** Write a function `forbes()` which takes as input  $F$ , the weight of fat tissue at a given time (i.e. the function  $F(t)$  evaluated at a certain time), and returns Forbe's equation given in (5.5). Also write the function `energy_balance()` which takes as input  $F$ ,  $L$ ,  $PAL$ , and  $EI$  and returns the energy balance as given in (5.8). In `energy_balance()` we also have that  $F$  is the fat tissue weight at a given time, and  $L$  is the lean tissue weight at a given time.

Using `forbes()` and `energy_balance()`, define the function `weight_odesystem()` which takes as input the current time as a float  $t$  and the current fat and lean weights as an array  $y$  and returns the right hand side of (5.4) as a tuple.

Use  $\rho_F = 9400$ ,  $\rho_L = 1800$ ,  $\gamma_F = 3.2$ ,  $\gamma_L = 22$ ,  $\eta_F = 180$ ,  $\eta_L = 230$ ,  $K = 0$  and  $\beta_{AT} = 0.14$ .

Hint: The functions `forbes()` and `energy_balance()` are not time dependent in the same way equations (5.5) and (5.8) are. The time dependent portions of these functions,  $F(t)$  and  $L(t)$ , are determined by what will be input from the  $y$  argument of `weight_odesystem()`.

**Problem 6.** Consider the initial value problem corresponding to (5.4). The following function returns the fat mass of an individual based on body weight (kg), age (years), height (meters), and sex. Use this function to define initial conditions  $F_0$  and  $L_0$  for the IVP above:  $F_0 = \text{fat\_mass(args*)}$ ,  $L_0 = BW - F_0$ .

```
def fat_mass(BW, age, H, sex):
    BMI = BW / H**2.
    if sex == 'male':
        return BW * (-103.91 + 37.31 * log(BMI) + 0.14 * age) / 100
    else:
        return BW * (-102.01 + 39.96 * log(BMI) + 0.14 * age) / 100
```

Suppose a 38 year old female, standing 5'8" and weighing 160 lbs, reduces her intake from 2143 to 2025 calories/day, and increases her physical activity from little to no exercise ( $PAL=1.4$ ) to exercising to 2-3 days per week ( $PAL=1.5$ ).

Use (5.9) and the original intake and physical activity levels to compute  $K$  for this system. Then use `solve_ivp` to solve the IVP. Graph the solution curve for this single-stage weightloss intervention over a period of 5 years. Your plot should match figure 5.3.

Note the provided code requires quantities in metric units (kilograms, meters, days) while our graph is converted to units of pounds and days. Use the conversions 1 lb = 2.204 kg, 1 ft = 0.305 m, and 1 yr = 365 days.

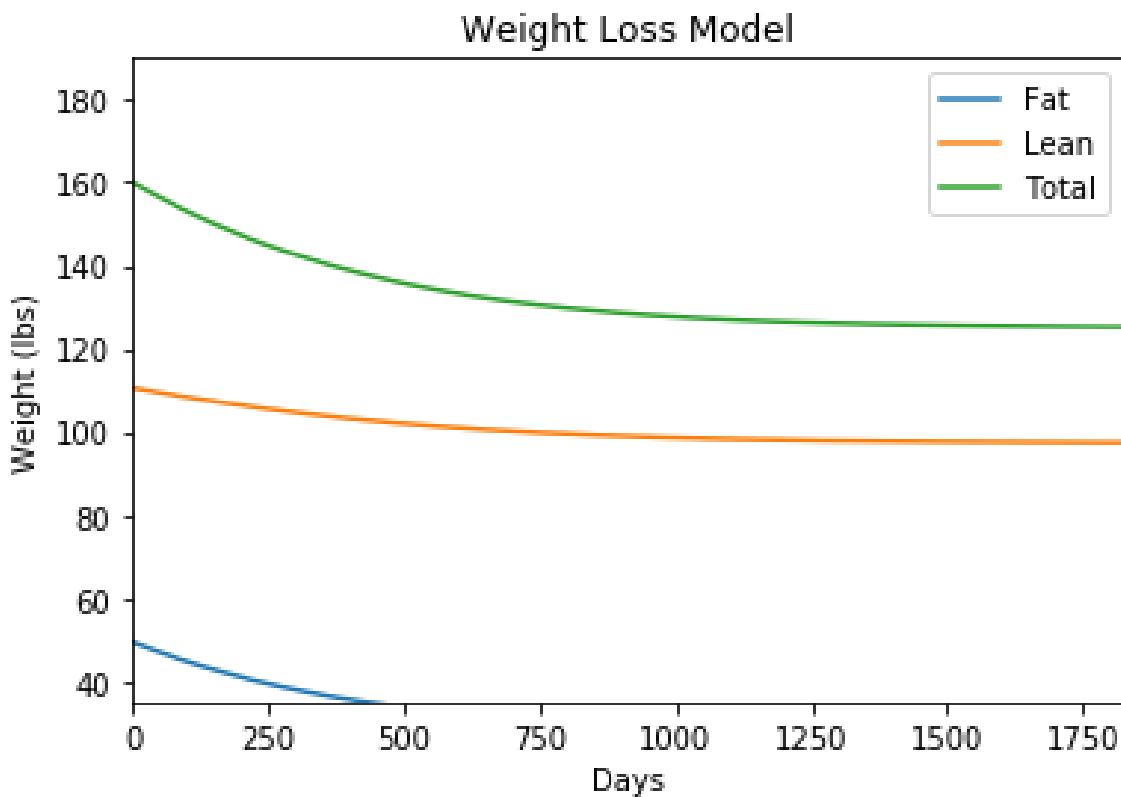


Figure 5.3: The solution of the weight change model for problem 6.

**Problem 7.** Modify the preceding problem to handle a two stage weightloss intervention: Suppose for the first 16 weeks intake is reduced from 2143 to 1600 calories/day and physical activity is increased from little to no exercise ( $PAL=1.4$ ) to an hour of exercise 5 days per week ( $PAL=1.7$ ). The following 16 weeks intake is increased from 1600 to 2025 calories/day, and exercise is limited to only 2-3 days per week ( $PAL=1.5$ ).

You will need to recompute  $F_0$ , and  $L_0$  at the end of the first 16 weeks, but  $K$  will stay the same. Find and graph the solution curve over the 32 week period.

# 6

# Lorenz Equations

**Lab Objective:** Investigate the behavior of a system that exhibits chaotic behavior. Demonstrate methods for visualizing the evolution of a system.

Chaos is everywhere. It can crop up in unexpected places and in remarkably simple systems, and a great deal of work has been done to describe the behavior of chaotic systems. One primary characterization of chaos is that small changes in initial conditions result in large changes over time in the solution curves.

## The Lorenz System

One of the earlier examples of chaotic behavior was discovered by Edward Lorenz. In 1963, while working to study atmospheric dynamics, he derived the simple system of equations

$$\begin{aligned}\frac{dx}{dt} &= \sigma(y - x) \\ \frac{dy}{dt} &= \rho x - y - xz \\ \frac{dz}{dt} &= xy - \beta z\end{aligned}$$

where  $\sigma$ ,  $\rho$ , and  $\beta$  are all constants. After deriving these equations, he plotted the solutions and observed some unexpected behavior. For appropriately chosen values of  $\sigma$ ,  $\rho$ , and  $\beta$ , the solutions did not tend toward any steady fixed points, nor did the system permit any stable cycles. The solutions did not tend off toward infinity either. This began the study of what was called a strange attractor. Although relatively simple, this system exhibits chaotic behavior.

**Problem 1.** Write a function that implements the Lorenz equations. Let  $\sigma = 10$ ,  $\rho = 28$ ,  $\beta = \frac{8}{3}$ . Make a 3D plot of a solution to the Lorenz equations, where the initial conditions  $x_0, y_0, z_0$  are each drawn randomly from a uniform distribution on  $[-15, 15]$  and for  $t$  in the range  $[0, 25]$ . As usual, use `scipy.integrate.solve_ivp` to compute an approximate solution. Compare your results with Figure 6.1.

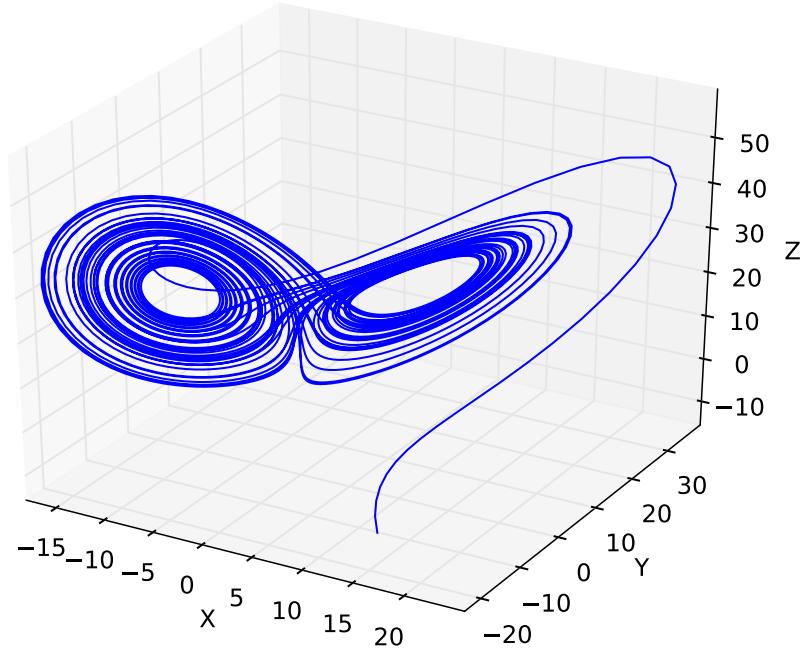


Figure 6.1: Approximate solution to the Lorenz equation with random initial conditions

## Basin of Attraction

Notice in the first problem that the solution tended to a certain region, called an attractor. The basin of attraction of an attractor is the set of initial conditions that tend towards the attractor. We will investigate the basin of attraction of the Lorenz system by changing the initial conditions.

**Problem 2.** To better visualize the Lorenz attractor, produce a single 3D plot displaying three solutions to the Lorenz equations, each with random initial conditions as before. Compare your results with Figure 6.2.

## Chaos

Chaotic systems exhibit high sensitivity to initial conditions. This means that a small difference in initial conditions will generally result in solutions that diverge significantly from each other. However, chaotic systems are not random. An explanation given by Lorenz is that chaos is "when the present determines the future, but the approximate present does not approximately determine the future."

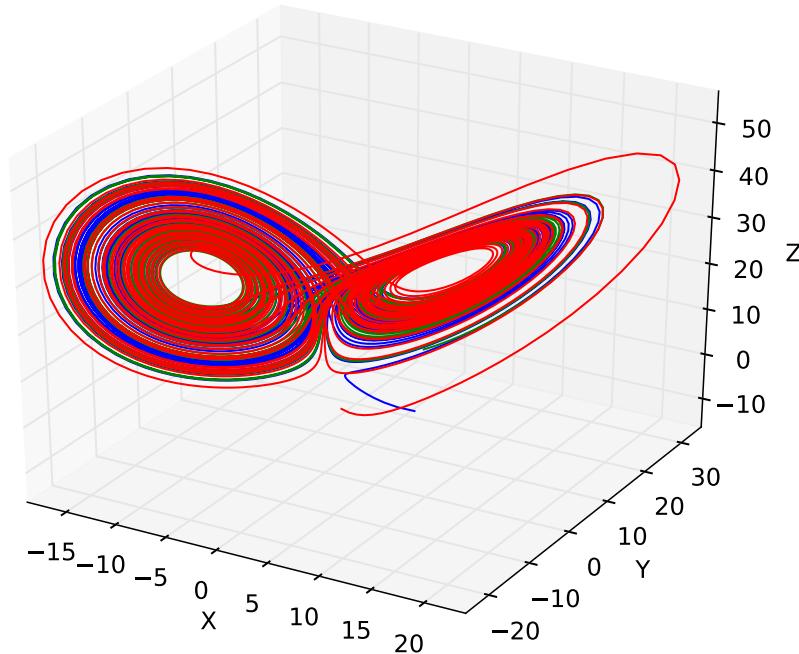


Figure 6.2: Multiple solutions to the Lorenz equation with random initial conditions

**Problem 3.** Use `matplotlib.animation.FuncAnimation` to produce a 3D animation of two solutions to the Lorenz equations with nearly identical initial conditions. To make the initial conditions, draw  $x_0, y_0, z_0$  as before, and then make a second initial condition by adding a small perturbation to the first (try using `np.random.randn(3)*(1e-10)` for the perturbation). Note that it may take several seconds before the separation between the two solutions will be noticeable.

The animation should display a point marker as well as the past trajectory curve for each solution. Save your animation as `lorenz_animation1.mp4`.

Hint: to ensure that the two numerical solutions have the same timesteps, you can use the `t_eval` argument of `solve_ivp`. Using this argument causes `solve_ivp` to return the solution's values at the points you pass in.

In a chaotic system, round-off error implicit in a numerical method can also cause divergent solutions. For example, using a Runge-Kutta method with two different values for the stepsize  $h$  on identical initial conditions will still result in approximations that differ in a chaotic fashion.

**Problem 4.** Even differences due to small numerical errors can cause solutions of chaotic systems to diverge from each other. The `solve_ivp` function allows user to specify error tolerances (similar to setting a value of  $h$  in a Runge-Kutta method). Using a single initial condition, produce two approximations by using the `solve_ivp` arguments (`atol=1e-15, rtol=1e-13`) for the first approximation and (`atol=1e-12, rtol=1e-10`) for the second.

As in the previous problem, use `FuncAnimation` to animation both solutions. Save the animation as `lorenz_animation2.mp4`.

## Lyapunov Exponents

The Lyapunov exponent of a dynamical system is one measure of how chaotic a system is. While there are more conditions for a system to be considered chaotic, one of the primary indicators of a chaotic system is extreme sensitivity to initial conditions. Strictly speaking, this is saying that a chaotic system is poorly conditioned. In a chaotic system, the sensitivity to changes in initial conditions depends exponentially on the time the system is allowed to evolve. If  $\delta(t)$  represents the difference between two solution curves, when  $\delta(t)$  is small, the following approximation holds.

$$\|\delta(t)\| \sim \|\delta(0)\| e^{\lambda t}$$

where  $\lambda$  is a constant called the Lyapunov exponent. In other words,  $\log(\|\delta(t)\|)$  is approximately linear as a function of time, with slope  $\lambda$ . For the Lorenz system (and for the parameter values specified in this lab), experimentally it can be verified that  $\lambda \approx .9$ .

**Problem 5.** Estimate the Lyapunov exponent of the Lorenz equations by doing the following:

- Produce an initial condition that already lies in the attractor. This can be done by using a random "dummy" initial condition, approximating the resulting solution to the Lorenz system for a short time, and then using the endpoint of that solution (which is now in the attractor) as the desired initial condition.
- Produce a second initial condition by adding a small perturbation to the first (as before).
- For both initial conditions, use `solve_ivp` to produce approximate solutions for  $0 \leq t \leq 10$
- Compute  $\|\delta(t)\|$  by taking the norm of the vector difference between the two solutions for each value of  $t$ .
- Use `scipy.stats.linregress` to calculate a best-fit line for  $\log(\|\delta(t)\|)$  against  $t$ .
- The slope of the best-fit line is an approximation for the Lyapunov exponent  $\lambda$

Produce a plot similar to Figure 6.3 by using `plt.semilogy`.

Hint: Remember that the best-fit line you calculated corresponds to a best-fit exponential for  $\|\delta(t)\|$ . If  $a$  and  $b$  are the slope and intercept of the best-fit line, the best-fit exponential can be plotted using `plt.semilogy(t, np.exp(a*t+b))`.

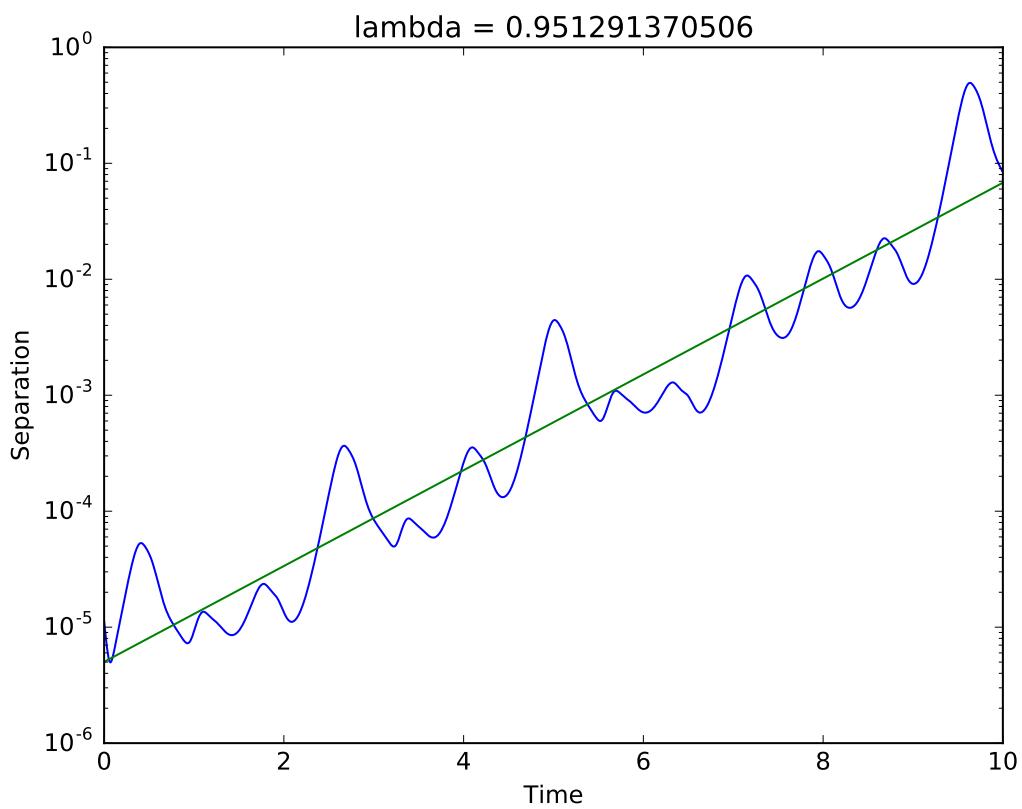


Figure 6.3: A semilog plot of the separation between two solutions to the Lorenz equations together with a fitted line that gives an estimate of the Lyapunov exponent of the system.



## 7

# Bifurcations and Hysteresis

Recall that any ordinary differential equation can be written as a first order system of ODEs,

$$x' = F(x), \quad x' := \frac{d}{dt}x(t). \quad (7.1)$$

Many interesting applications and physical phenomena can be modeled using ODEs. Given a mathematical model of the form (7.1), it is important to understand geometrically how its solutions behave. This information can then be conveyed in a phase portrait, a graph describing solutions of (7.1) with differential initial conditions. The first step in constructing a phase portrait is to find the equilibrium solutions of the equation, i.e., the zeros of  $F(x)$ , and to determine their stability.

It is often the case that the mathematical model we study depends on some parameter or set of parameters  $\lambda$ . Thus the ODE becomes

$$x' = F(x, \lambda). \quad (7.2)$$

The parameter  $\lambda$  can then be tuned to better fit the physical application. As  $\lambda$  varies, the equilibrium solutions and other geometric features of (7.2) may suddenly change. A value of  $\lambda$  where the phase portrait changes is called a bifurcation point; the study of how these changes occur is called bifurcation theory. The parameter values and corresponding equilibrium solutions are often graphed together in a bifurcation diagram.

As an example, consider the scalar differential equation

$$x' = x^2 + \lambda. \quad (7.3)$$

For  $\lambda > 0$  equation (7.3) has no equilibrium solutions. At  $\lambda = 0$  the equilibrium point  $x = 0$  appears, and for  $\lambda < 0$  it splits into two equilibrium points. For this system, a bifurcation occurs at  $\lambda = 0$ . This is an example of a saddle-node bifurcation. The bifurcation diagram is shown in Figure 7.1

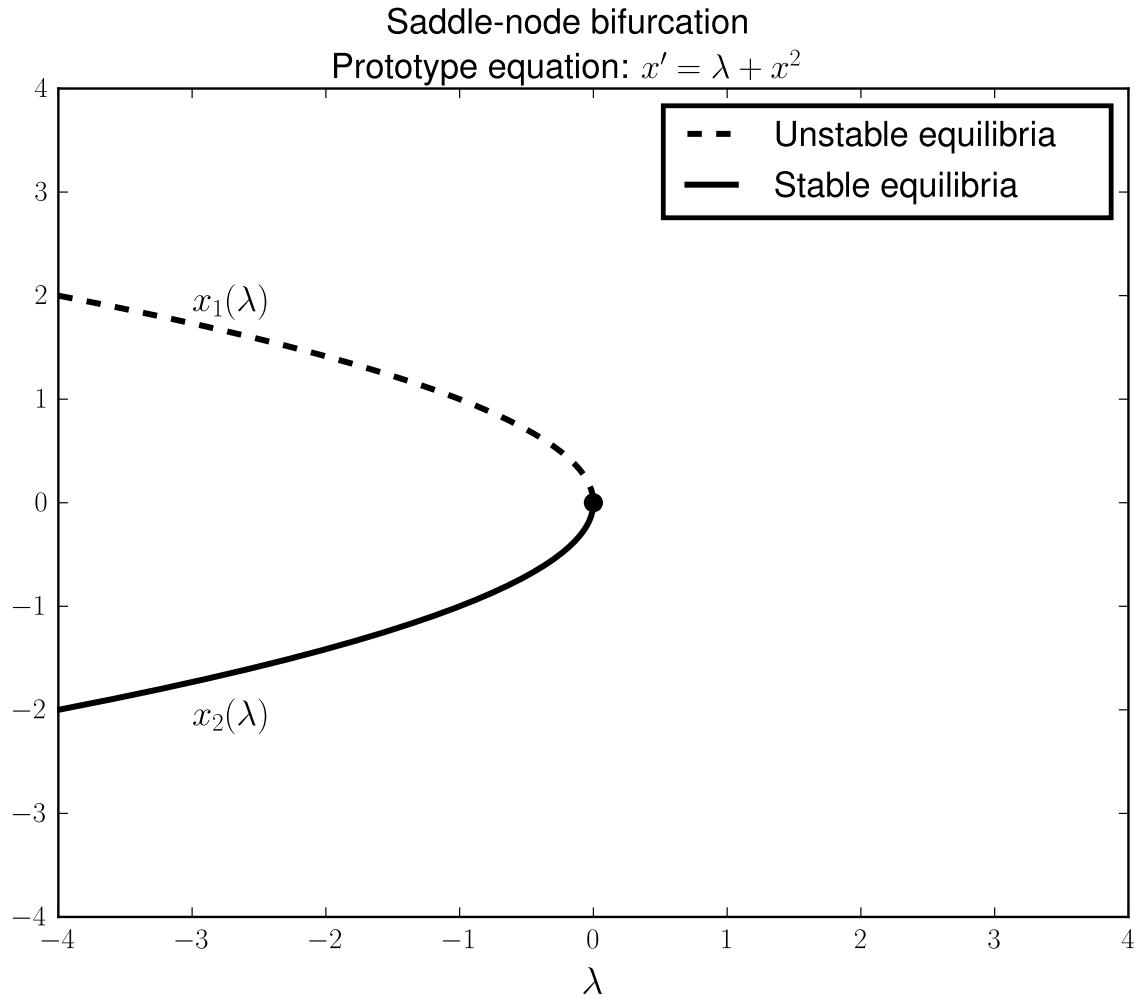


Figure 7.1: Bifurcation diagram for the equation  $x' = \lambda + x^2$ .

Suppose that  $F(x_0, \lambda_0) = 0$ . We use a method called natural embedding to find zeros  $(x, \lambda)$  of  $F$  for nearby values of  $\lambda$ . Specifically, we step forward in  $\lambda$  by letting  $\lambda_1 = \lambda_0 + \Delta\lambda$ , and use Newton's method to find the value  $x_1$  that satisfies  $F(x_1, \lambda_1) = 0$ . This method works well except when  $\lambda$  is near a bifurcation point  $\lambda^*$ .

The following code implements the natural embedding algorithm, and then uses that algorithm to find the curves in the bifurcation diagram for (7.3). Notice that this algorithm needs a good initial guess for  $x_0$  to get started.

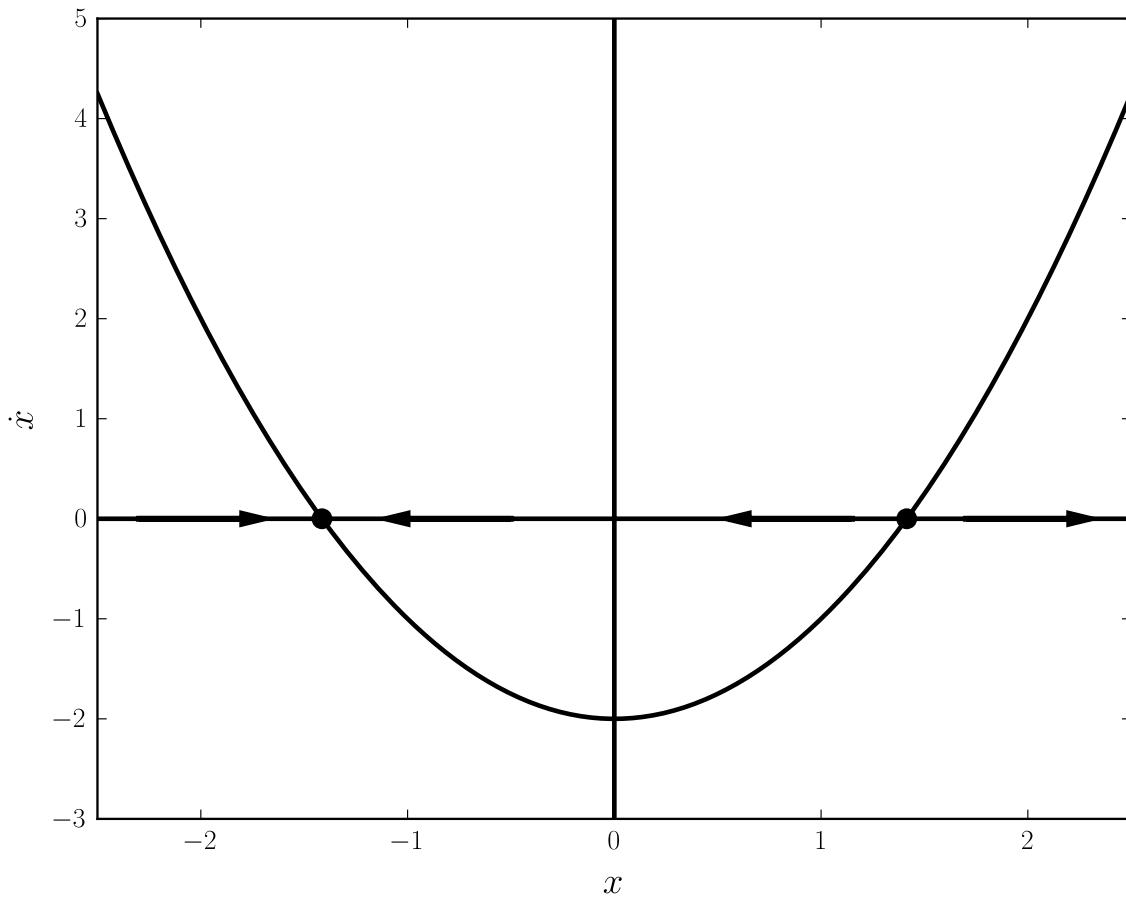


Figure 7.2: Phase Portrait for the equation  $x' = -2 + x^2$ .

```

import numpy as np
import matplotlib.pyplot as plt
from scipy.optimize import newton

def EmbeddingAlg(param_list, guess, F):
    X = []
    for param in param_list:
        try:
            # Solve for x_value making F(x_value, param) = 0.
            x_value = newton(F, guess, fprime=None, args=(param,), tol=1E-7, ←
                maxiter=50)
            # Record the solution and update guess for the next iteration.
            X.append(x_value)
            guess = x_value
        except RuntimeError:
            # If Newton's method fails, return a truncated list of parameters
            # with the corresponding x values.

```

```

        return param_list[:len(X)], X
    # Return the list of parameters and the corresponding x values.
    return param_list, X

def F(x, lmbda):
    return x**2 + lmbda

# Top curve shown in the bifurcation diagram
C1, X1 = EmbeddingAlg(np.linspace(-5, 0, 200), np.sqrt(5), F)
# The bottom curve
C2, X2 = EmbeddingAlg(np.linspace(-5, 0, 200), -np.sqrt(5), F)

```

**Problem 1.** Use the natural embedding algorithm to create a bifurcation diagram for the differential equation

$$x' = \lambda x - x^3.$$

This type of bifurcation is called a pitchfork bifurcation (you should see a pitchfork in your diagram).

Hints: Essentially this amounts to running the same code as the example, but with different parameters and function calls so that you are tracing through the right curves for this problem. To make this first problem work, you will want to have your ‘linspace’ run from high to low instead of from low to high. There will be three different lines in this image all of which must be produced using the EmbeddingAlg function. Any hard coding will result in an automatic 0. See Figure 7.3.

**Problem 2.** Create bifurcation diagrams for the differential equation

$$x' = \eta + \lambda x - x^3,$$

where  $\eta = -1, -0.2, 0.2$  and  $1$ . Notice that when  $\eta = 0$  you can see the pitchfork bifurcation of the previous problem. There should be four different images, one for each value of  $\eta$ . Each image will be built of 3 pieces. See Figure 7.4.

## Hysteresis

The following ODE exhibits an interesting bifurcation phenomenon called hysteresis:

$$x' = \lambda + x - x^3.$$

This system has a bifurcation diagram containing what is known as a hysteresis loop, shown in Figure 7.5. In the hysteresis loop, when the parameter  $\lambda$  moves beyond the bifurcation point the equilibrium solution makes a sudden jump to the other stable branch. When this occurs the system cannot reach its previous equilibrium by simply rewinding the parameter slightly. The next section discusses a model with a hysteresis loop.

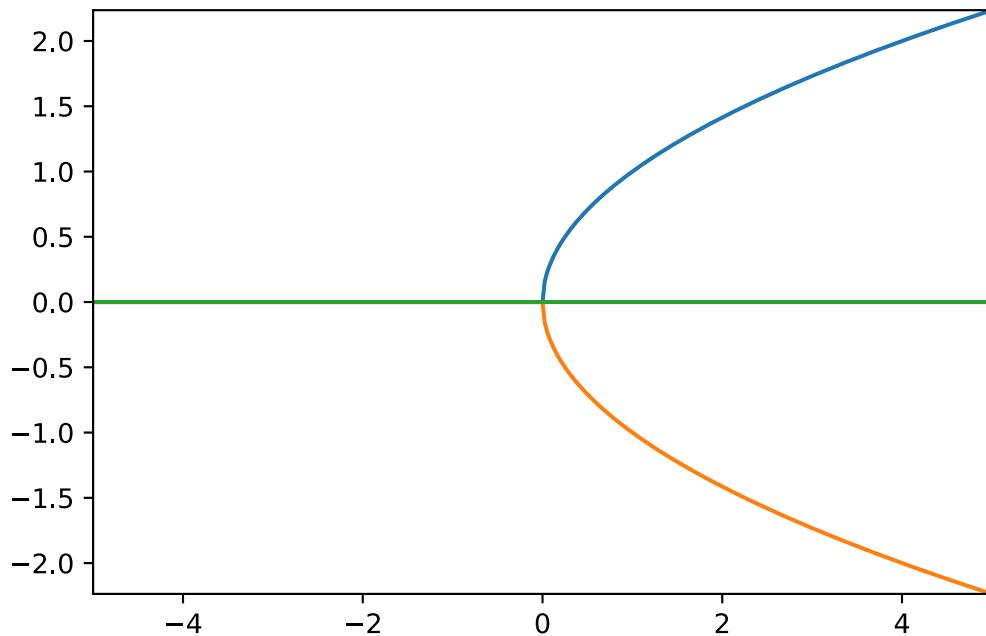


Figure 7.3: Bifurcation diagram for the equation  $x' = \lambda x - x^3$ .

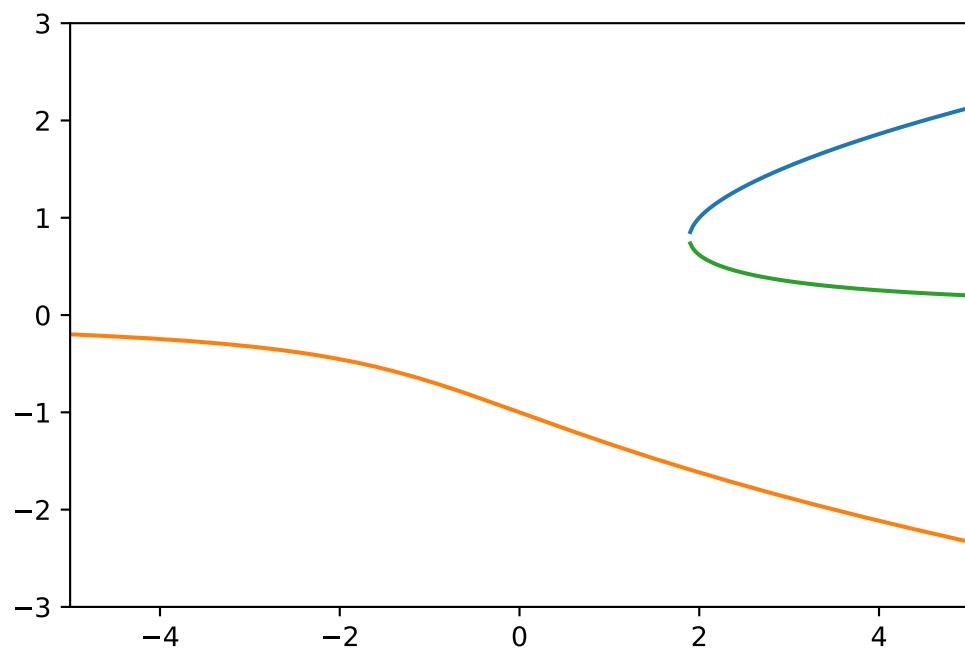


Figure 7.4: Bifurcation diagram for the equation  $x' = \eta + \lambda x - x^3$  with  $\eta = -1$ .

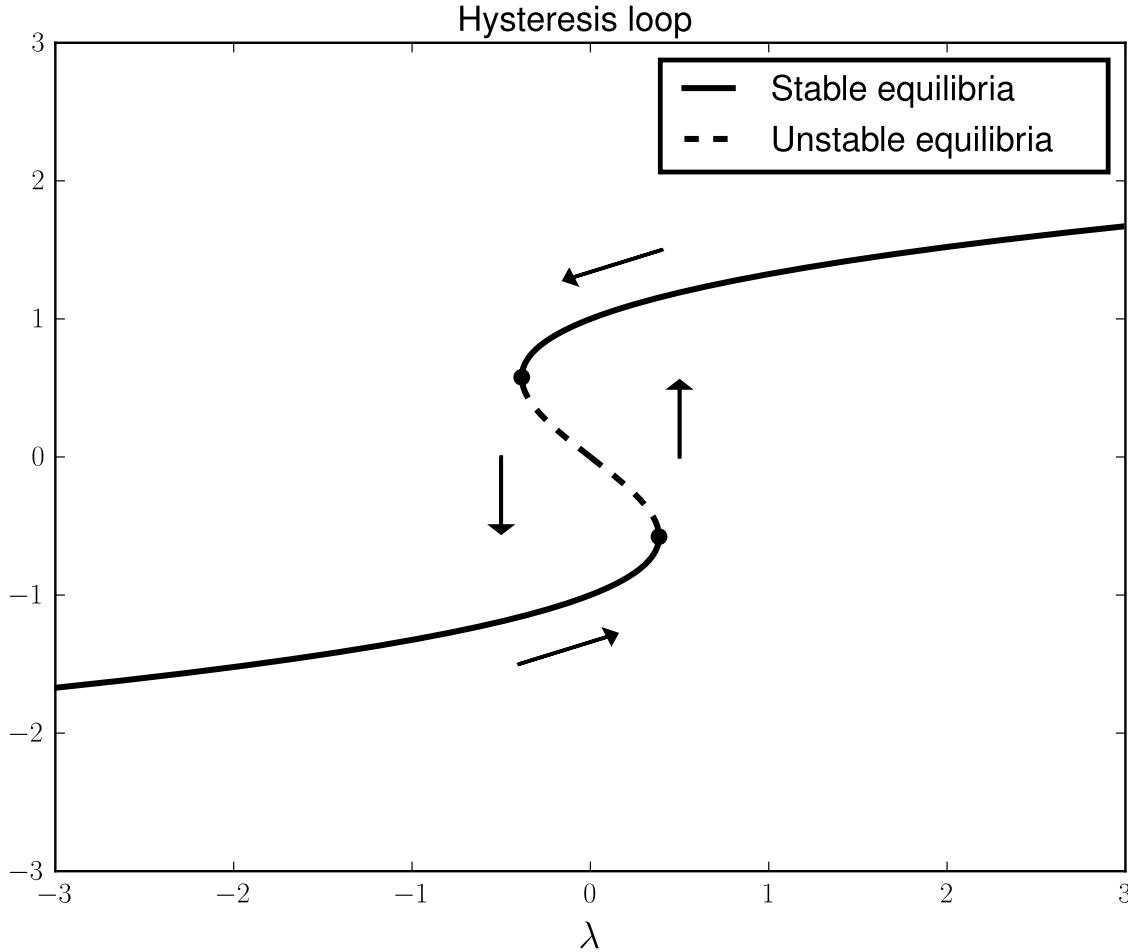


Figure 7.5: Bifurcation diagram for the ODE  $x' = \lambda + x - x^3$ .

### Budworm Population Dynamics

Here we study a mathematical model describing the population dynamics of an insect called the spruce budworm. In eastern Canada, an outbreak in the budworm population can destroy most of the trees in a forest of balsam fir trees in about 4 years. The mathematical model is given by

$$N' = RN \left(1 - \frac{N}{K}\right) - p(N). \quad (7.4)$$

This model was studied by Ludwig et al (1978), and is described well in Strogatz's text Nonlinear Dynamics and Chaos. Here  $N(t)$  represents the budworm population at time  $t$ ,  $R$  is the growth rate of the budworm population and  $K$  represents the carrying capacity of the environment. We could interpret  $K$  to represent the amount of food available to the budworms.  $p(N)$  represents the death rate of budworms due to predators (birds); we assume specifically that  $p(N)$  has the form  $P(N) = \frac{BN^2}{A^2+N^2}$ .

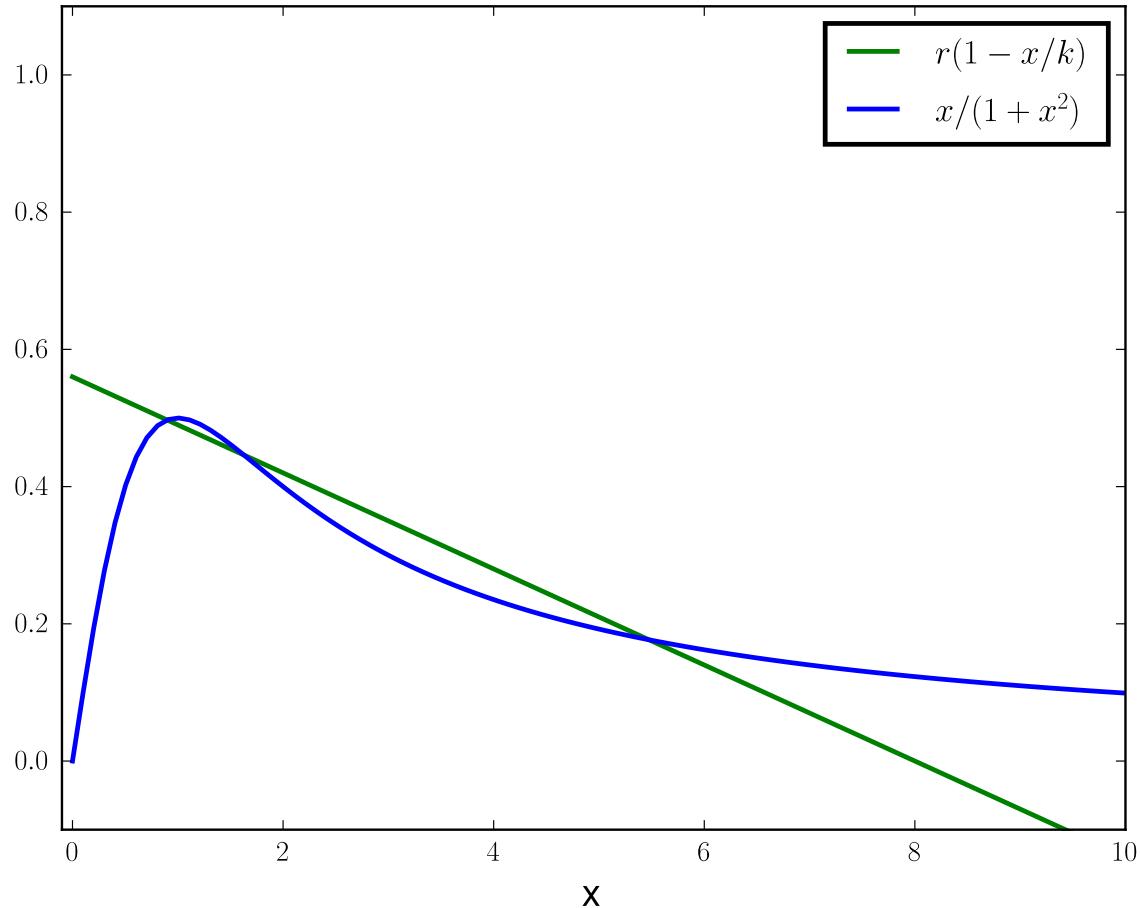


Figure 7.6: Graphical demonstration of nonzero equilibrium solutions for the budworm population (here  $r = .56$ ,  $k = 8$ ); equilibrium solutions occur where the curves cross. As  $k$  increases, the line  $y = r(1 - x/k)$  gets more shallow and the number of solutions goes from one to three and then back to one.

Before studying the equilibrium points of (7.4) it is important to reduce the number of parameters in the system by nondimensionalizing. Thus, we make the coordinate change  $x = N/A$ ,  $\tau = Bt/A$ ,  $r = RA/B$ , and  $k = K/A$ , obtaining finally the system

$$\frac{dx}{d\tau} = rx(1 - x/k) - \frac{x^2}{1 + x^2}. \quad (7.5)$$

Note that  $x = 0$  is always an equilibrium solution. To find other equilibrium solutions we study the equation  $r(1 - x/k) - x/(1 + x^2) = 0$ . Fix  $r = .56$ , and consider Figure (7.6) ( $k = 8$  in the figure).

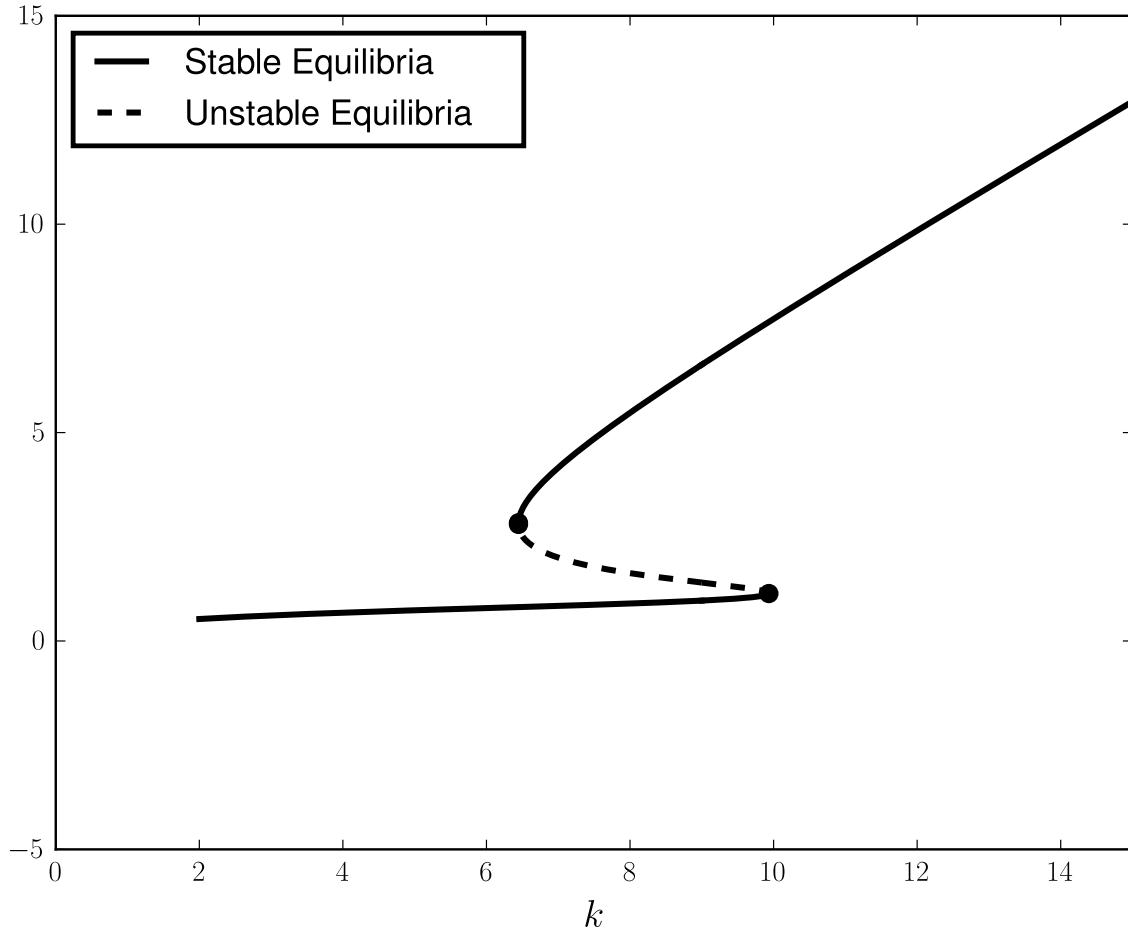


Figure 7.7: Bifurcation diagram for the budworm population model. The parameter  $r$  is fixed at 0.56. The lower stable branch is known as the refuge level of the budworm population, while the upper stable branch is known as the outbreak level. Once the budworm population reaches an outbreak level, the available food (foliage of the balsam fir trees) in the system must be reduced drastically to jump back down to refuge level. Thus many of the balsam fir trees die before the budworm population returns to refuge level.

**Problem 3 (Budworm Population).** Reproduce the bifurcation diagram (7.5) for the differential equation

$$\frac{dx}{d\tau} = rx(1 - x/k) - \frac{x^2}{1 + x^2},$$

where  $r = 0.56$ .

Hint: Find a value for  $k$  that you know is in the middle of the plot (i.e. where there are three possible solutions), then use the code from the previous problems to expand along each contour till you obtain the desired curve. Now find the proper initial guesses that give you the right bifurcation curve. The final plot will look like the one in Figure 7.7, but you will have to run the embedding algorithm 4-6 times to get every part of the plot. In order to make a black dashed line, add ' $-k$ ' as the third argument in `plt.plot()` and use ' $-k$ ' as the third argument in `plt.plot()` to get the solid black line.

# 8

## The Finite Difference Method

**Lab Objective:** The finite difference method provides a solid foundation for solving partial differential equations. Understanding and applying finite difference is key to understanding numerical solutions to PDEs.

A **finite difference** for a function  $f(x)$  is an expression of the form  $f(x + s) - f(x + t)$ . Finite differences can give a good approximation of derivatives.

Suppose we have a function  $u(x)$ , defined on an interval  $[a, b]$ . Let  $a = x_0, x_1, \dots, x_{n-1}, x_n = b$  be a grid of  $n + 1$  evenly spaced points, with  $x_{i+1} - x_i = h$ , where  $h = (b - a)/n$ .

You are used to seeing the derivative  $u'(x)$ , which can written in centered-difference form as:

$$u'(x) = \lim_{h \rightarrow 0} \frac{u(x + h) - u(x - h)}{2h}.$$

Suppose we are interested in knowing the value of the derivative at the points  $\{x_i\}$ . Even if we don't have a formula for  $u'(x)$ , we can approximate it using finite differences. We first write the Taylor polynomial expansion of  $u(x + h)$  and  $u(x - h)$  centered at  $x$ . This gives

$$u(x + h) = u(x) + u'(x)h + \frac{1}{2}u''(x)h^2 + \frac{1}{6}u'''(x)h^3 + \mathcal{O}(h^4) \quad (8.1)$$

$$u(x - h) = u(x) - u'(x)h + \frac{1}{2}u''(x)h^2 - \frac{1}{6}u'''(x)h^3 + \mathcal{O}(h^4) \quad (8.2)$$

Subtracting (8.2) from (8.1) and rearranging gives

$$u'(x) = \frac{u(x + h) - u(x - h)}{2h} + \mathcal{O}(h^2).$$

In terms of our grid points  $\{x_i\}$ , we have:

$$u'(x_i) \approx \frac{u(x_i + h) - u(x_i - h)}{2h} = \frac{u(x_{i+1}) - u(x_{i-1})}{2h}.$$

We won't worry about the derivative at the endpoints,  $u'(x_0)$  and  $u'(x_n)$ . This allows us to approximate the values  $\{u'(x_i)\}$  as the solution to a system of equations:

$$\frac{1}{2h} \begin{bmatrix} -1 & 0 & 1 & & \\ & -1 & 0 & 1 & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 0 & 1 \\ & & & -1 & 0 & 1 \end{bmatrix}_{(n-1) \times (n+1)} \cdot \begin{bmatrix} u(x_0) \\ u(x_1) \\ \vdots \\ u(x_{n-1}) \\ u(x_n) \end{bmatrix}_{(n+1) \times 1} \approx \begin{bmatrix} u'(x_1) \\ u'(x_2) \\ \vdots \\ u'(x_{n-2}) \\ u'(x_{n-1}) \end{bmatrix}_{(n-1) \times 1}. \quad (8.3)$$

This can be rewritten with a  $(n-1) \times (n-1)$  tridiagonal matrix instead:

$$\frac{1}{2h} \begin{bmatrix} 0 & 1 & & \\ -1 & 0 & 1 & \\ & \ddots & \ddots & \ddots \\ & & -1 & 0 & 1 \\ & & & -1 & 0 \end{bmatrix}_{(n-1) \times (n-1)} \cdot \begin{bmatrix} u(x_1) \\ u(x_2) \\ \vdots \\ u(x_{n-2}) \\ u(x_{n-1}) \end{bmatrix}_{(n-1) \times 1} + \begin{bmatrix} -u(x_0)/(2h) \\ 0 \\ \vdots \\ 0 \\ u(x_n)/(2h) \end{bmatrix}_{(n-1) \times 1} \approx \begin{bmatrix} u'(x_1) \\ u'(x_2) \\ \vdots \\ u'(x_{n-2}) \\ u'(x_{n-1}) \end{bmatrix}_{(n-1) \times 1}. \quad (8.4)$$

Next, we will consider the approximation for  $u''(x)$ . If we let

$$u'(x) \approx \frac{u(x + \frac{h}{2}) - u(x - \frac{h}{2})}{h}$$

then

$$\begin{aligned} u''(x) &\approx \frac{u'(x + \frac{h}{2}) - u'(x - \frac{h}{2})}{h} \approx \frac{\frac{u((x + \frac{h}{2}) + \frac{h}{2}) - u((x + \frac{h}{2}) - \frac{h}{2})}{h} - \frac{u((x - \frac{h}{2}) + \frac{h}{2}) - u((x - \frac{h}{2}) - \frac{h}{2})}{h}}{h} \\ &= \frac{u(x + h) - 2u(x) + u(x - h)}{h^2} \end{aligned}$$

You can achieve the same result by again consider the Taylor polynomial expansion and adding (8.1) and (8.2) and rearranging. Thus

$$u''(x_i) \approx \frac{u(x_i + h) - 2u(x_i) + u(x_i - h)}{h^2} = \frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1})}{h^2}$$

Again ignoring the second derivative at the endpoints, this can be written in matrix form as

$$\frac{1}{h^2} \begin{bmatrix} 1 & -2 & 1 & & \\ & 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots & \\ & & 1 & -2 & 1 \\ & & & 1 & -2 & 1 \end{bmatrix}_{(n-1) \times (n+1)} \cdot \begin{bmatrix} u(x_0) \\ u(x_1) \\ \vdots \\ u(x_{n-1}) \\ u(x_n) \end{bmatrix}_{(n+1) \times 1} \approx \begin{bmatrix} u''(x_1) \\ u''(x_2) \\ \vdots \\ u''(x_{n-2}) \\ u''(x_{n-1}) \end{bmatrix}_{(n-1) \times 1}. \quad (8.5)$$

This can also be written with a  $(n-1) \times (n-1)$  tridiagonal matrix:

$$\frac{1}{h^2} \begin{bmatrix} -2 & 1 & & \\ 1 & -2 & 1 & \\ & \ddots & \ddots & \ddots \\ & & 1 & -2 & 1 \\ & & & 1 & -2 \end{bmatrix}_{(n-1) \times (n-1)} \cdot \begin{bmatrix} u(x_1) \\ u(x_2) \\ \vdots \\ u(x_{n-2}) \\ u(x_{n-1}) \end{bmatrix}_{(n-1) \times 1} + \begin{bmatrix} u(x_0)/h^2 \\ 0 \\ \vdots \\ 0 \\ u(x_n)/h^2 \end{bmatrix}_{(n-1) \times 1} = \begin{bmatrix} u''(x_1) \\ u''(x_2) \\ \vdots \\ u''(x_{n-2}) \\ u''(x_{n-1}) \end{bmatrix}_{(n-1) \times 1}. \quad (8.6)$$

**Problem 1.** Let  $u(x) = \sin((x + \pi)^2 - 1)$ . Use (8.3) - (8.6) to approximate  $\frac{1}{2}u'' - u'$  at the grid points where  $a = 0$ ,  $b = 1$ , and  $n = 10$ . Graph the result.

The previous equations are not only useful for approximating derivatives, but they can be also used to solve differential equations. Suppose that instead of knowing the function  $u(x)$ , we know that  $\frac{1}{2}u'' - u' = f$ , where the function  $f(x)$  is given. How do we solve for  $u(x)$ ?

## Finite Difference Methods

Numerical methods for differential equations seek to approximate the exact solution  $u(x)$  at some finite collection of points in the domain of the problem. Instead of analytically solving the original differential equation, defined over an infinite-dimensional function space, they use a well-chosen finite system of algebraic equations to approximate the original problem.

Consider the following differential equation:

$$\begin{aligned} \varepsilon u''(x) - u'(x) &= f(x), \quad x \in (0, 1), \\ u(0) &= \alpha, \quad u(1) = \beta. \end{aligned} \tag{8.7}$$

Equation (8.7) can be written  $Du = f$ , where  $D = \varepsilon \frac{d^2}{dx^2} - \frac{d}{dx}$  is a differential operator defined on the infinite-dimensional space of functions that are twice continuously differentiable on  $[0, 1]$  and satisfy  $u(0) = \alpha$ ,  $u(1) = \beta$ .

We look for an approximate solution  $\{U_i\}$ , where

$$U_i \approx u(x_i)$$

on an evenly spaced grid of points,  $a = x_0, x_1, \dots, x_n = b$ . Our finite difference method will replace the differential operator  $D = \varepsilon \frac{d^2}{dx^2} - \frac{d}{dx}$ , (which is defined on an infinite-dimensional space), with finite difference operators (defined on a finite dimensional space). To do this, we replace derivative terms in the differential equation with appropriate difference expressions.

Recalling that

$$\begin{aligned} \frac{d^2}{dx^2}u(x_i) &= \frac{u(x_{i+1}) - 2u(x_i) + u(x_{i-1})}{h^2} + \mathcal{O}(h^2), \\ \frac{d}{dx}u(x_i) &= \frac{u(x_{i+1}) - u(x_{i-1})}{2h} + \mathcal{O}(h^2). \end{aligned}$$

we define the finite difference operator  $D_h$  by

$$D_h U_i = \frac{1}{h^2} (U_{i+1} - 2U_i + U_{i-1}) - \frac{1}{2h} (U_{i+1} - U_{i-1}). \tag{8.8}$$

Thus we discretize equation (8.7) using the equations

$$\frac{\varepsilon}{h^2} (U_{i+1} - 2U_i + U_{i-1}) - \frac{1}{2h} (U_{i+1} - U_{i-1}) = f(x_i), \quad i = 1, \dots, n-1,$$

along with boundary conditions  $U_0 = \alpha$ ,  $U_n = \beta$ .

This gives  $n + 1$  equations and  $n + 1$  unknowns, and can be written in matrix form as

$$\frac{1}{h^2} \begin{bmatrix} h^2 & 0 & 0 & \dots & 0 \\ (\varepsilon + h/2) & -2\varepsilon & (\varepsilon - h/2) & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & \dots & (\varepsilon + h/2) & -2\varepsilon & (\varepsilon - h/2) \\ 0 & \dots & & 0 & h^2 \end{bmatrix}_{(n+1) \times (n+1)} \cdot \begin{bmatrix} U_0 \\ U_1 \\ \vdots \\ U_{n-1} \\ U_n \end{bmatrix}_{(n+1) \times 1} = \begin{bmatrix} \alpha \\ f(x_1) \\ \vdots \\ f(x_{n-1}) \\ \beta \end{bmatrix}_{(n+1) \times 1}.$$

As before, we can remove two equations to modify the system to obtain an  $(n-1) \times (n-1)$  tridiagonal system:

$$\frac{1}{h^2} \begin{bmatrix} -2\varepsilon & (\varepsilon - h/2) & 0 & \dots & 0 \\ (\varepsilon + h/2) & -2\varepsilon & (\varepsilon - h/2) & \dots & 0 \\ \vdots & & \ddots & & \vdots \\ 0 & \dots & (\varepsilon + h/2) & -2\varepsilon & (\varepsilon - h/2) \\ 0 & \dots & & (\varepsilon + h/2) & -2\varepsilon \end{bmatrix}_{(n-1) \times (n-1)} \cdot \begin{bmatrix} U_1 \\ U_2 \\ \vdots \\ U_{n-2} \\ U_{n-1} \end{bmatrix}_{(n-1) \times 1} = \begin{bmatrix} f(x_1) - \alpha(\varepsilon + h/2)/h^2 \\ f(x_2) \\ \vdots \\ f(x_{n-2}) \\ f(x_{n-1}) - \beta(\varepsilon - h/2)/h^2 \end{bmatrix}_{(n-1) \times 1}. \quad (8.9)$$

**Problem 2.** Use equation (8.9) to solve the singularly perturbed BVP (8.7) on the interval  $[0, 1]$  with  $\varepsilon = 1/10$ ,  $f(x) = -1$ ,  $\alpha = 1$ , and  $\beta = 3$  on a grid with  $n = 30$  subintervals. Graph the solution. This BVP is called singularly perturbed because of the location of the parameter  $\varepsilon$ . For  $\varepsilon = 0$  the ODE has a drastically different character - it then becomes first order, and can no longer support two boundary conditions.

## A heuristic test for convergence

The finite differences used above are second order approximations of the first and second derivatives of a function. It seems reasonable to expect that the numerical solution would converge at a rate of about  $\mathcal{O}(h^2)$ . How can we check that a numerical approximation is reasonable?

Suppose a finite difference method is  $\mathcal{O}(h^p)$  accurate. This means that the error  $E(h) \approx Ch^p$  for some constant  $C$  as  $h \rightarrow 0$  (in other words, for  $h > 0$  small enough).

So compute the approximation  $y_k$  for each stepsize  $h_k$ ,  $h_1 > h_2 > \dots > h_m$ .  $y_m$  should be the most accurate approximation, and will be thought of as the true solution. Then the error of the approximation for stepsize  $h_k$ ,  $k < m$ , is

$$E(h_k) = \max(|y_k - y_m|) \approx Ch_k^p,$$

$$\log(E(h_k)) = \log(C) + p \log(h_k).$$

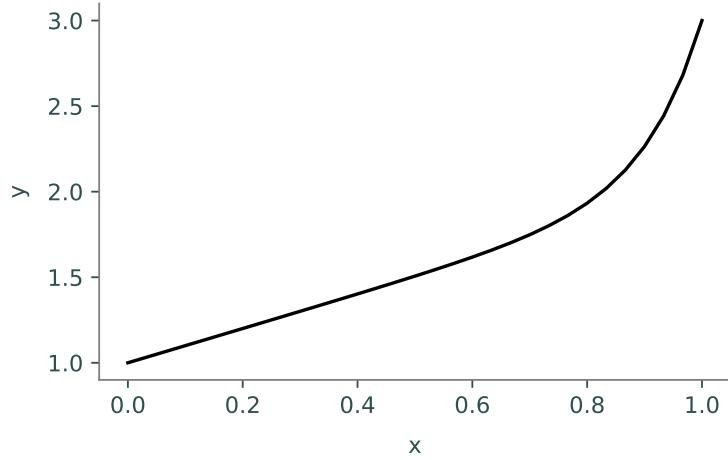


Figure 8.1: The solution to Problem 2. The solution gets steeper near  $x = 1$  as  $\varepsilon$  gets small.

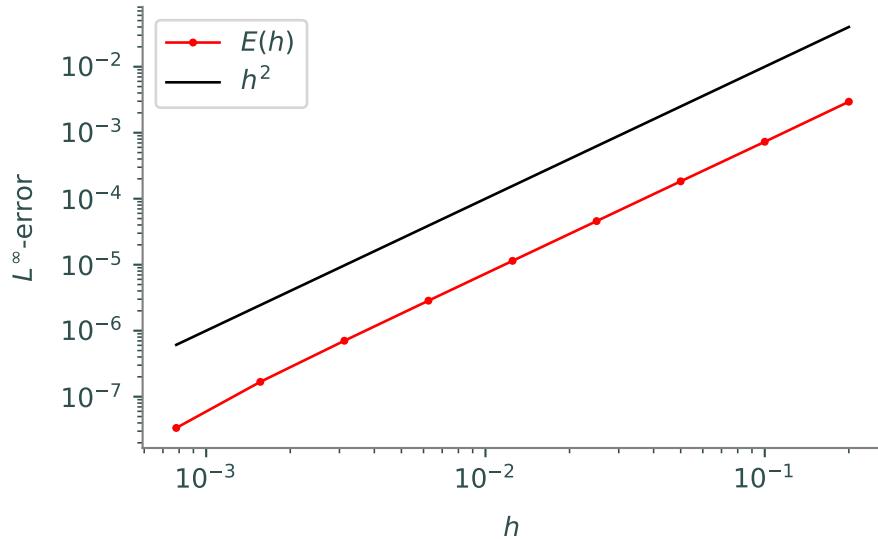


Figure 8.2: Demonstration of second order convergence for the finite difference approximation (8.8) of the BVP given in (8.7) with  $\varepsilon = .5$ .

Thus on a log-log plot of  $E(h)$  vs.  $h$ , these values should be on a straight line with slope  $p$  when  $h$  is small enough to start getting convergence. We should note that demonstrating second-order convergence does NOT imply that the numerical approximation is converging to the correct solution.

**Problem 3.** Implement a function `singular_bvp` to compute the finite difference solution to 8.7. Using  $n = 5 \times 2^0, 5 \times 2^1, \dots, 5 \times 2^9$  subintervals, compute 10 approximate solutions. Use these to visualize the  $\mathcal{O}(h^2)$  convergence of the finite difference method from Problem 2 by producing a loglog plot of error against subinterval count; this will be similar to Figure 8.2, except with  $\varepsilon = 0.1$ .

To produce the plot, treat the approximation with  $n = 5 \times 2^9$  subintervals as the "true solution", and measure the error for the other approximations against it. Note that, since the ratios of numbers of subintervals between approximations are multiples of 2, we can compute the  $L_\infty$  error for the  $n = 5 \times 2^j$  approximation by using the `step` argument in the array slicing syntax:

```
# best approximation; the vector has length 5*2^9+1
sol_best = singular_bvp(eps, alpha, beta, f, 5*(2**9))

# approximation with 5*(2^j) intervals; the vector has length 5*2^j+1
sol_approx = singular_bvp(eps, alpha, beta, f, 5*(2**j))

# approximation error; slicing results in a vector of length 5*2^j+1,
# which allows it to be compared
error = np.max(np.abs(sol_approx - sol_best[::2**(9-j)]))
```

**Problem 4.** Extend your finite difference code to the case of a general second order linear BVP with boundary conditions:

$$\begin{aligned} a_1(x)y''(x) + a_2(x)y'(x) + a_3(x)y(x) &= f(x), \quad x \in (a, b), \\ y(a) &= \alpha, \quad y(b) = \beta. \end{aligned}$$

Use your code to solve the boundary value problem

$$\begin{aligned} \varepsilon y'' - 4(\pi - x^2)y &= \cos x, \\ y(0) &= 0, \quad y(\pi/2) = 1, \end{aligned}$$

for  $\varepsilon = 0.1$  on a grid with  $n = 30$  subintervals. Be sure to modify the finite difference operator  $D_h$  in (8.8) correctly.

The next few problems will help you test your finite difference code.

**Problem 5.** Numerically solve the boundary value problem

$$\begin{aligned} \varepsilon y''(x) + xy'(x) &= -\varepsilon\pi^2 \cos(\pi x) - \pi x \sin(\pi x), \\ y(-1) &= -2, \quad y(1) = 0, \end{aligned}$$

for  $\varepsilon = 0.1, 0.01$ , and  $0.001$ . Use a grid with  $n = 150$  subintervals. Plot your solutions.

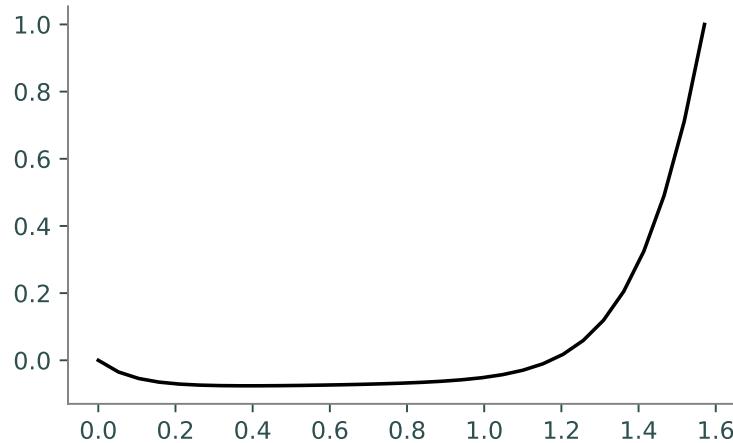


Figure 8.3: The solution to Problem 4.

**Problem 6.** Numerically solve the boundary value problem

$$(\varepsilon + x^2)y''(x) + 4xy'(x) + 2y(x) = 0, \\ y(-1) = 1/(1 + \varepsilon), \quad y(1) = 1/(1 + \varepsilon),$$

for  $\varepsilon = 0.05, 0.02$ . Use a grid with  $n = 150$  subintervals. Plot your solutions.

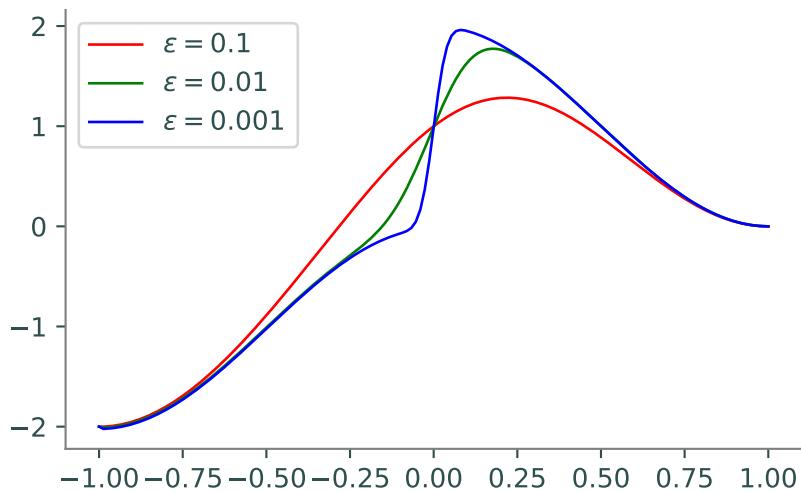


Figure 8.4: The solution to Problem 5.

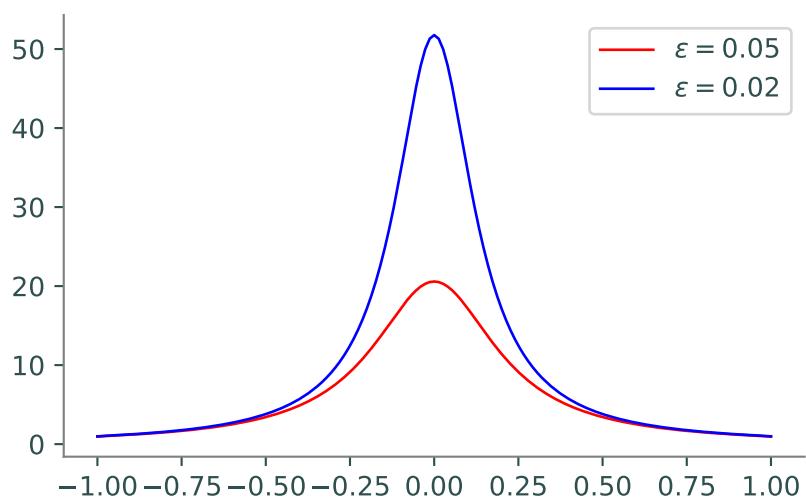


Figure 8.5: The solution to Problem 6.

# 9

# Wave Phenomena

## Advection Equation

The advection equation (or transport equation) is given by  $u_t + su_x = 0$ , where  $s$  is a nonzero constant. Consider the Cauchy problem

$$\begin{aligned} u_t + su_x &= 0, \quad -\infty < x < \infty, \\ u(x, 0) &= f(x). \end{aligned}$$

The function  $f(x)$  may be thought of as an initial wave or signal. The general solution of this initial boundary value problem is  $u(x, t) = f(x - st)$  (check this!). The solution  $u(x, t)$  is a traveling wave that takes the signal  $f(x)$  and moves it along at a constant speed  $s$  — to the right if  $s > 0$ , and to the left if  $s < 0$ .

## Wave Equation

Many different wave phenomena can be described using a hyperbolic PDE called the wave equation. These wave phenomena occur in fields such as electromagnetics, fluid dynamics, and acoustics. This equation is given by

$$u_{tt} = s^2 \Delta u. \tag{9.1}$$

The 1D equation can be derived in the context of many physical models; a common derivation describes the motion of a string vibrating in a plane. Another nice derivation uses Hooke's law from the theory of elasticity.

After making the change of variables  $(\xi, \eta) = (x - st, x + st)$  and using the chain rule, we find that the 1D wave equation  $u_{tt} = s^2 u_{xx}$  is equivalent to  $u_{\xi\eta} = 0$ . The general solution of this last equation is

$$u(\xi, \eta) = F(\xi) + G(\eta)$$

for some scalar functions  $F$  and  $G$ . In  $(x, t)$  coordinates the solution is

$$u(x, t) = F(x - st) + G(x + st)$$

Thus the general solution of the wave equation is the sum of two parts: one is a signal traveling to the right with constant speed  $|s|$ , and the other is a signal traveling to the left with speed  $|s|$ .

The wave equation is usually seen in the context of an initial boundary value problem. This takes the form

$$\begin{aligned} u_{tt} &= s^2 u_{xx}, \quad 0 < x < l, \quad t > 0, \\ u(0, t) &= u(l, t) = 0, \\ u(x, 0) &= f(x), \\ u_t(x, 0) &= g(x). \end{aligned}$$

### Numerical solution of the wave equation

We look to approximate  $u(x, t)$  on a grid of points  $(x_j, t_m)_{j=0, m=0}^{J, M}$ . Denote the approximation to  $u(x_j, t_m)$  by  $U_j^m$ . Recall that the centered approximations in space and time are

$$\begin{aligned} D_{tt} U_j^m &= \frac{U_j^{m+1} - 2U_j^m + U_j^{m-1}}{(\Delta t)^2}, \\ D_{xx} U_j^m &= \frac{U_{j+1}^m - 2U_j^m + U_{j-1}^m}{(\Delta x)^2}. \end{aligned}$$

The resulting method is given by

$$\begin{aligned} \frac{U_j^{m+1} - 2U_j^m + U_j^{m-1}}{(\Delta t)^2} &= s^2 \frac{U_{j+1}^m - 2U_j^m + U_{j-1}^m}{(\Delta x)^2}, \\ U_j^{m+1} &= -U_j^{m-1} + 2(1 - \lambda^2)U_j^m + \lambda^2(U_{j+1}^m + U_{j-1}^m), \end{aligned}$$

where  $\lambda = s(\Delta t)/(\Delta x)$ . This method may be written in matrix form as

$$U^{m+1} = AU^m - U^{m-1}$$

where

$$A = \begin{bmatrix} 2(1 - \lambda^2) & \lambda^2 & & \\ \lambda^2 & 2(1 - \lambda^2) & \lambda^2 & \\ & \ddots & \ddots & \ddots \\ & & \lambda^2 & 2(1 - \lambda^2) & \lambda^2 \\ & & & \lambda^2 & 2(1 - \lambda^2) \end{bmatrix}$$

and

$$U^m = \begin{bmatrix} U_1^m \\ U_2^m \\ \vdots \\ U_{J-1}^m \end{bmatrix}$$

In the matrix equation above, we have already used the boundary conditions to determine that  $U_0^m = U_J^m = 0$  at each time  $t_m$ . Note that, to obtain the approximation  $U_j^{m+1}$  of  $u(x_j, t_{m+1})$ , the method uses the value of the approximation at the previous two time steps. We can find the solution for the first two time steps by using the initial conditions. Using the initial conditions directly gives an approximation at  $t = t_0 = 0$ :

$$U_j^0 = f(x_j), \quad 1 \leq j \leq J-1$$

To obtain an approximation at the second time step, we consider the Taylor expansion

$$u(x_j, t_1) = u(x_j, 0) + u_t(x_j, 0)\Delta t + u_{tt}(x_j, 0)\frac{\Delta t^2}{2} + u_{ttt}(x_j, t_1^*)\frac{\Delta t^3}{6}.$$

Recalling that the solution  $u(x, t)$  satisfies the wave equation, we substitute in expressions from our initial conditions:

$$u(x_j, t_1) = u(x_j, 0) + g(x_j)\Delta t + s^2 f''(x_j)\frac{\Delta t^2}{2} + u_{ttt}(x_j, t_1^*)\frac{\Delta t^3}{6}.$$

Ignoring the third order term, we obtain a second order approximation for the second time step:

$$U_j^1 = U_j^0 + g(x_j)\Delta t + s^2 f''(x_j)\frac{\Delta t^2}{2}, \quad 1 \leq j \leq J-1$$

or if  $f$  is not readily differentiable,

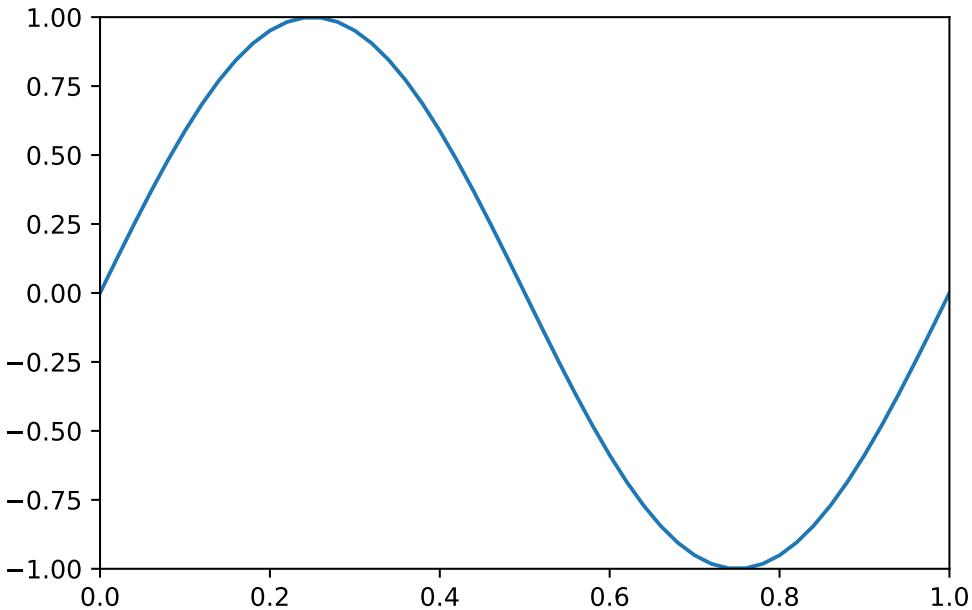
$$U_j^1 = U_j^0 + g(x_j)\Delta t + \frac{\lambda^2}{2}(U_{j-1}^0 - 2U_j^0 + U_{j+1}^0)$$

This method is conditionally stable; the CFL condition (a necessary condition for convergence of PDEs) is that  $\lambda \leq 1$ .

**Problem 1.** Consider the initial boundary value problem

$$\begin{aligned} u_{tt} &= u_{xx}, \\ u(0, t) &= u(1, t) = 0, \\ u(x, 0) &= \sin(2\pi x), \\ u_t(x, 0) &= 0. \end{aligned}$$

Numerically approximate the solution  $u(x, t)$  for  $t \in [0, .5]$ . Use  $J = 50$  subintervals in the  $x$  dimension and  $M = 50$  subintervals in the  $t$  dimension. Animate the results. Compare your results with the analytic solution  $u(x, t) = \sin(2\pi x) \cos(2\pi t)$  graphically. This function is known as a standing wave. See Figure 9.1.

Figure 9.1:  $u(x, t = 0)$ .

**Problem 2.** Consider the initial boundary value problem

$$\begin{aligned} u_{tt} &= u_{xx}, \\ u(0, t) &= u(1, t) = 0, \\ u(x, 0) &= .2e^{-m^2(x-1/2)^2} \\ u_t(x, 0) &= .4m^2(x - 1/2)e^{-m^2(x-1/2)^2}. \end{aligned}$$

The solution of this problem is a Gaussian pulse. It travels to the right at a constant speed. This solution models, for example, a wave pulse in a stretched string. Note that the fixed boundary conditions reflect the pulse back when it meets the boundary.

Numerically approximate the solution  $u(x, t)$  for  $t \in [0, 1]$ . Set  $m = 20$ . Use 200 subintervals in space and 220 in time, and animate your results. Then use 200 subintervals in space and 180 in time, and animate your results. Note that the stability condition is not satisfied for the second mesh. See 9.2.

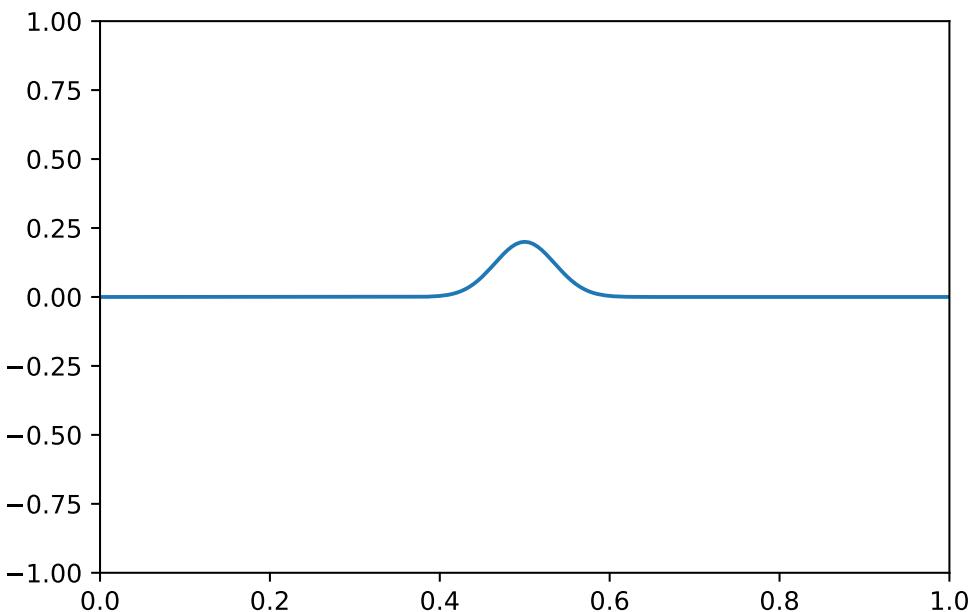


Figure 9.2:  $u(x, t = 0)$ .

**Problem 3.** Consider the initial boundary value problem

$$\begin{aligned} u_{tt} &= u_{xx}, \\ u(0, t) &= u(1, t) = 0, \\ u(x, 0) &= .2e^{-m^2(x-1/2)^2} \\ u_t(x, 0) &= 0. \end{aligned}$$

The initial condition separates into two smaller, slower-moving pulses, one traveling to the right and the other to the left. This solution models, for example, a plucked guitar string

Numerically approximate the solution  $u(x, t)$  for  $t \in [0, 2]$ . Set  $m = 20$ . Use 200 subintervals in space and 440 in time, and animate your results. It is rather easy to see that the solution to this problem is the sum of two traveling waves, one traveling to the left and the other to the right, as described earlier.

**Problem 4.** Consider the initial boundary value problem

$$\begin{aligned} u_{tt} &= u_{xx}, \\ u(0, t) &= u(1, t) = 0, \\ u(x, 0) &= \begin{cases} 1/3 & \text{if } 5/11 < x < 6/11, \\ 0 & \text{otherwise} \end{cases} \\ u_t(x, 0) &= 0. \end{aligned}$$

Numerically approximate the solution  $u(x, t)$  for  $t \in [0, 2]$ . Use 200 subintervals in space and 440 in time, and animate your results. Even though the method is second order and stable for this discretization, since the initial condition is discontinuous there are large dispersive errors. See Figure 9.3.

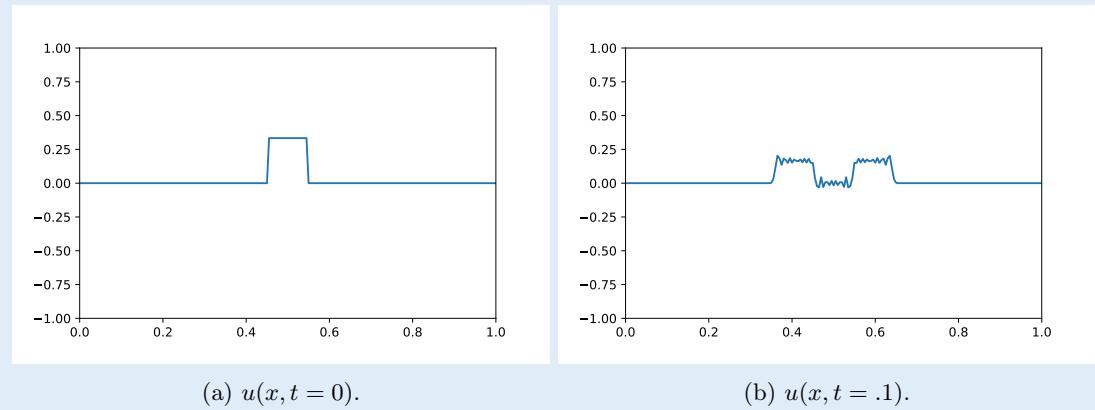


Figure 9.3: The graphs for Problem 4 at various times  $t$ .

## Traveling Wave Solutions of an Evolution Equation

Recall that the advection (transport) equation with initial conditions, given by

$$\begin{aligned} u_t + su_x &= 0, \quad -\infty < x < \infty, \\ u(x, 0) &= f(x), \end{aligned}$$

has as its general solution  $u(x, t) = f(x - st)$ . Consider a general evolutionary PDE of the form

$$u_t = G(u, u_x, u_{xx}, \dots) \tag{9.2}$$

An interesting question to ask is whether (9.2) has traveling wave solutions: is there a signal or wave profile  $f(x)$ , so that  $u(x, t) = f(x - st)$  is a solution of (9.2) that carries the signal at a constant speed  $s$ ? These traveling waves are often significant physically. For example, in a PDE modeling insect population dynamics a traveling wave could represent a swarm of locusts; in a PDE describing a combustion process a traveling wave could represent an explosion or detonation.

## Burgers' equation

We will examine the process of studying traveling wave solutions using Burgers' equation, a nonlinear PDE from gas dynamics. It is given by

$$u_t + \left( \frac{u^2}{2} \right)_x = \nu u_{xx}, \quad (9.3)$$

where  $u$  and  $\nu$  represent the velocity and viscosity of the gas, respectively. It models both the process of transport with the nonlinear advection term  $(u^2/2)_x = uu_x$ , as well as diffusion due to the viscosity of the gas  $(\nu u_{xx})$ .

Let us look for a traveling wave solution  $u(x, t) = \hat{u}(x - st)$  for Burgers equation. We transform (9.3) into the moving frame  $(x, t) \rightarrow (\bar{x}, \bar{t}) = (x - st, t)$ . In this frame (9.3) becomes

$$u_{\bar{t}} - su_{\bar{x}} + \left( \frac{u^2}{2} \right)_{\bar{x}} = \nu u_{\bar{x}\bar{x}} \quad (9.4)$$

This new frame of reference corresponds to an observer moving along with the wave, so that the wave appears stationary as the observer studies it. Thus,  $\hat{u}_{\bar{t}} = 0$ , so that the wave profile  $\hat{u}$  satisfies the ordinary differential equation

$$-su_{\bar{x}} + \left( \frac{u^2}{2} \right)_{\bar{x}} = \nu u_{\bar{x}\bar{x}}. \quad (9.5)$$

From here on we will drop the bar notation for simplicity. We seek a traveling wave solution with asymptotically constant boundary conditions; that is,  $\lim_{x \rightarrow \pm\infty} \hat{u}(x) = u_{\pm}$  both exist, and  $\lim_{x \rightarrow \pm\infty} \hat{u}'(x) = 0$ . We will suppose that  $u_- > u_+ > 0$ .

Note that to this point we still don't know the speed of the traveling wave. Integrating both sides of this differential equation, and then taking the limit as  $x \rightarrow +\infty$ , we obtain

$$\begin{aligned} -s \int_{-\infty}^x u' + \int_{-\infty}^x \left( \frac{u^2}{2} \right)' &= \nu \int_{-\infty}^x u'', \\ -s(u(x) - u_-) + \frac{u^2(x)}{2} - \frac{u_-^2}{2} &= \nu(u'(x) - u'(-\infty)), \\ -s(u_+ - u_-) + \frac{u_+^2}{2} - \frac{u_-^2}{2} &= 0. \end{aligned}$$

Thus given boundary conditions  $u_{\pm}$  at  $\pm\infty$ , the speed of the traveling wave must be  $s = \frac{u_- + u_+}{2}$ .

Usually at this point, the traveling wave must be numerically solved using the profile ODE ((9.5) for Burgers equation). However, the profile ODE for Burgers is simple enough that it is possible to obtain an analytic solution. The traveling wave is given by

$$\hat{u}(x) = s - a \tanh \left( \frac{ax}{2\nu} + \delta \right)$$

where  $a = (u_- - u_+)/2$  and  $\delta$  is fixed real number. We get a family of solutions because any translation of a traveling wave solution is also a traveling wave solution.

## Stability of traveling waves

Suppose that an evolutionary PDE

$$u_t = G(u, u_x, u_{xx}, \dots). \quad (9.6)$$

has a traveling wave solution  $u(x, t) = \hat{u}(x - st)$ . An interesting question to consider is whether the mathematical solution,  $\hat{u}$ , has a physical analogue. In other words, does the traveling wave show up in real life? This question is the start of the mathematical study of stability of traveling waves.

We begin by translating (9.6) into the moving frame  $(x, t) \rightarrow (\bar{x}, \bar{t}) = (x - st, t)$ . In this frame the PDE becomes

$$u_t - su_x = G(u, u_x, u_{xx}, \dots).$$

In these coordinates the traveling wave is stationary. Thus, the solution of

$$\begin{aligned} u_t - su_x &= G(u, u_x, u_{xx}, \dots), \\ u(x, t=0) &= \hat{u}(x), \end{aligned}$$

is given by  $u(x, t) = \hat{u}(x)$ . We say that the traveling wave  $\hat{u}$  is asymptotically orbitally stable if whenever  $v(x)$  is a small perturbation of  $\hat{u}(x)$ , the general solution of

$$\begin{aligned} u_t - su_x &= G(u, u_x, u_{xx}, \dots), \\ u(x, t=0) &= v(x), \end{aligned}$$

converges to some translation of  $\hat{u}$  as  $t \rightarrow \infty$ . Using this definition to prove stability of a traveling wave is a nontrivial task.

### Visualizing stability of the traveling wave solution of Burgers' equation

The traveling wave solution of Burgers' equation is a stable wave. To view this numerically, we discretize the PDE

$$u_t - su_x + uu_x = u_{xx}$$

using the second order centered approximations

$$\begin{aligned} D_t U_j^{n+1/2} &= \frac{U_j^{n+1} - U_j^n}{\Delta t}, \quad D_{xx} U_j^{n+1/2} = \frac{1}{2} \left( \frac{U_{j+1}^{n+1} - U_{j-1}^{n+1}}{2\Delta x} + \frac{U_{j+1}^n - U_{j-1}^n}{2\Delta x} \right), \\ D_{xx} U_j^{n+1/2} &= \frac{1}{2} \left( \frac{U_{j+1}^{n+1} - U_j^{n+1} + U_{j-1}^{n+1}}{(\Delta x)^2} + \frac{U_{j+1}^n - U_j^n + U_{j-1}^n}{(\Delta x)^2} \right) \end{aligned}$$

Substituting these expressions into the PDE we obtain a second-order, implicit Crank-Nicolson method

$$\begin{aligned} U_j^{n+1} - U_j^n &= K_1 [(s - U_j^{n+1})(U_{j+1}^{n+1} - U_{j-1}^{n+1}) + (s - U_j^n)(U_{j+1}^n - U_{j-1}^n)] \\ &\quad + K_2 [(U_{j+1}^{n+1} - 2U_j^{n+1} + U_{j-1}^{n+1}) + (U_{j+1}^n - 2U_j^n + U_{j-1}^n)], \end{aligned}$$

where  $K_1 = \frac{\Delta t}{4\Delta x}$  and  $K_2 = \frac{\Delta t}{2(\Delta x)^2}$ .

**Problem 5.** Numerically solve the initial value problem

$$\begin{aligned} u_t - su_x + uu_x &= u_{xx}, \quad x \in (-\infty, \infty), \\ u(x, 0) &= \hat{u}(x) + v(x), \end{aligned}$$

for  $t \in [0, 1]$ . Let the perturbation  $v(x)$  be given by

$$v(x) = 3.5(\sin(3x) + 1) \frac{1}{\sqrt{2\pi}} \exp(-x^2/2)$$

And let the initial condition be  $u(x, 0) = \hat{u}(x) + v(x)$ . Approximate the  $x$  domain,  $(-\infty, \infty)$ , numerically by the finite interval  $[-20, 20]$ , and fix  $\hat{u}(-20) = u_-$ ,  $\hat{u}(20) = u_+$ . Let  $u_- = 5$ ,  $u_+ = 1$  which makes  $s = 3$ . Use 150 intervals in space and 350 steps in time. Animate your results. You should see the solution converge to a translate of the traveling wave  $\hat{u}$ . See Figure 9.4.

Hint: This difference scheme is no longer a linear equation. We have a nonlinear equation in  $U^{n+1}$ . We can still solve this function using Newton's method or some other similar solver. In this case, use `scipy.optimize.fsolve`.

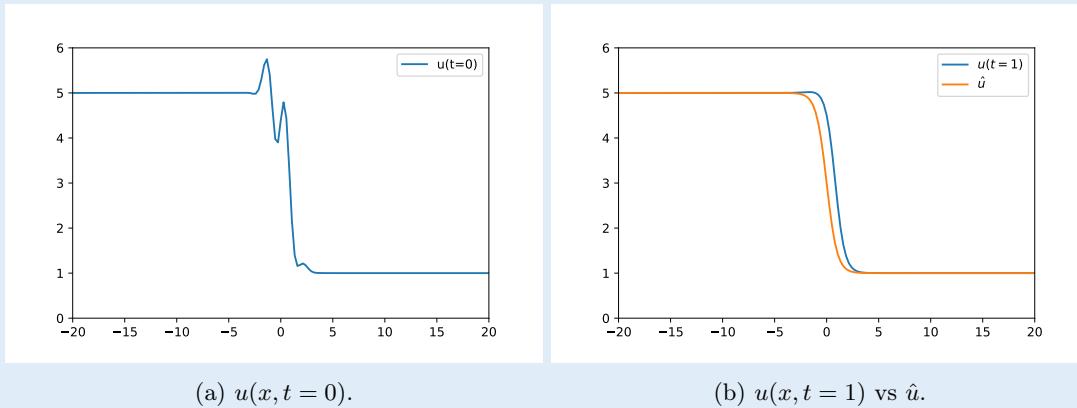


Figure 9.4: The graphs for Problem 5



# 10

## Conservation Laws and Heat Flow

Many physical phenomena have a conservation law associated with them. For instance, matter, energy, and momentum are all conserved quantities. The fundamental conservation law states that the rate of change of the total quantity in the system is equal to the rate that the quantity enters the system plus the rate at which the quantity is produced by sources inside the system. While this is a global property, we can use it to obtain a local differential equation that the concentration of the quantity must obey everywhere in the system. Because of this, conservation laws are very important in modeling a wide variety of phenomena.

### Derivation of the Conservation equation in multiple dimensions

Suppose  $\Omega$  is a region in  $\mathbb{R}^n$ , and  $V \subset \Omega$  is bounded with a reasonably well-behaved boundary  $\partial V$ . Let  $u(\vec{x}, t)$  represent the density (concentration) of some quantity throughout  $\Omega$ . Let  $\vec{n}(x)$  represent the normal direction to  $V$  at  $x \in \partial V$ , and let  $\vec{J}(\vec{x}, t)$  be the flux vector for the quantity, so that  $\vec{J}(\vec{x}, t) \cdot \vec{n}(x) dA$  represents the rate at which the quantity leaves  $V$  by crossing a boundary element with area  $dA$ . Note that the total amount of the quantity in  $V$  is

$$\int_V u(\vec{x}, t) dt,$$

and the rate at which the quantity enters  $V$  is

$$-\int_{\partial V} \vec{J}(\vec{x}, t) \cdot \vec{n}(x) dA.$$

We let the source term be given by  $g(\vec{x}, t, u)$ ; we may interpret this to mean that the rate at which the quantity is produced in  $V$  is

$$\int_V g(\vec{x}, t, u) dt.$$

Then the integral form of the conservation law for  $u$  is expressed as

$$\frac{d}{dt} \int_V u(\vec{x}, t) d\vec{x} = - \int_{\partial V} \vec{J} \cdot \vec{n} dA + \int_V g(\vec{x}, t, u) d\vec{x}.$$

If  $u$  and  $J$  are sufficiently smooth functions, then we have

$$\frac{d}{dt} \int_V u d\vec{x} = \int_V u_t d\vec{x},$$

and

$$\int_{\partial V} \vec{J} \cdot \vec{n} dA = \int_V \nabla \cdot \vec{J} d\vec{x}.$$

Putting these together yields

$$\int_V u(\vec{x}, t) d\vec{x} = \int_V (-\nabla \cdot \vec{J} + g(\vec{x}, t, u)) d\vec{x}$$

Since this holds for all nice subsets  $V \subset \Omega$  with  $V$  arbitrarily small, the integrands must be equal everywhere, and we obtain the differential form of the conservation law for  $u$ :

$$u_t + \nabla \cdot \vec{J} = g(\vec{x}, t, u),$$

where  $\nabla$  is the gradient operator and  $\nabla \cdot \vec{J} = \frac{\partial J_1}{\partial x_1} + \cdots + \frac{\partial J_n}{\partial x_n}$

## Constitutive Relations

So far, our conservation law consists of 2 unknowns ( $u$  and  $J$ ) but only 1 equation. To this equation we need to add other equations, called constitutive relations, which are used to fully determine the system.

For example, suppose we wish to model the flow of heat. Since heat flows from warmer regions to colder regions, and the rate of heat flow depends on the difference in temperature between regions, we usually assume that the flux vector  $\vec{J}$  is given by

$$\vec{J}(x, t) = -\nu \nabla u(x, t),$$

where  $\nu$  is called the diffusion constant and  $\nabla u(x, t) = [\partial_{x_1} u, \dots, \partial_{x_n} u]^T$ . This constitutive relation is called Fick's law, and is the basic model for any diffusive process. Substituting into the conservation law we obtain

$$u_t - \nu \Delta u(x, t) = g(\vec{x}, t, u)$$

where  $\Delta$  is the Laplacian operator:

$$\Delta u(x, t) = \frac{\partial^2 u}{\partial x_1^2} + \cdots + \frac{\partial^2 u}{\partial x_n^2}.$$

The function  $g$  represents heat sources and sinks within the region.

## Numerically modeling heat flow

Consider the heat flow equation in one dimension together with an appropriate initial condition  $u(x, 0) = f(x)$ , homogeneous Dirichlet boundary conditions, and  $g(x, t, u) = 0$ :

$$\begin{aligned} u_t &= \nu u_{xx}, \quad x \in [a, b], \quad t \in [0, T], \\ u(a, t) &= 0, \quad u(b, t) = 0, \\ u(x, 0) &= f(x). \end{aligned}$$

We will create an approximation  $U_i^j$  to  $u(x_i, t_j)$  on the grid  $x_i = a + hi$ ,  $t_j = kj$ , where  $h$  and  $k$  are small changes in  $x$  and  $t$  respectively and  $i$  and  $j$  are indices; so,  $U_i^j$  denotes the approximate value of  $u$  at the  $i$ -th grid point and the  $j$ -th time step.

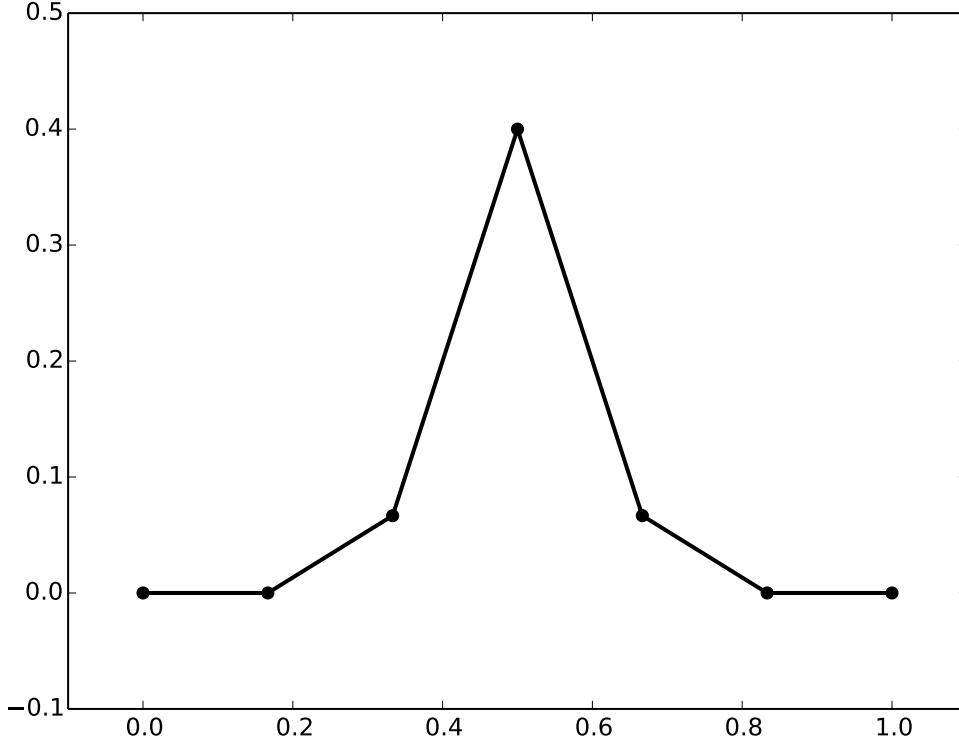


Figure 10.1: The graph of  $U^0$ , the approximation to the solution  $u(x, t = 0)$  for Problem 1.

As before, we will use the finite difference method to create this approximation. Recall that by using Taylor's theorem, we have the first-order forward difference approximation

$$u_t(x, t) = \frac{u(x, t + k) - u(x, t)}{k} + \mathcal{O}(k).$$

and the second-order centered difference approximation

$$u_{xx}(x_i, t_j) = \frac{u(x_i + h, t_j) - 2u(x_i, t_j) - u(x_i - h, t_j)}{h^2} + \mathcal{O}(h^2).$$

Applying these difference approximations give us the  $\mathcal{O}(h^2 + k)$  explicit method

$$\begin{aligned} \frac{U_i^{j+1} - U_i^j}{k} &= \nu \frac{U_{i+1}^j - 2U_i^j + U_{i-1}^j}{h^2}, \\ U_i^{j+1} &= U_i^j + \frac{\nu k}{h^2} (U_{i+1}^j - 2U_i^j + U_{i-1}^j). \end{aligned} \tag{10.1}$$

This method can be written in matrix form as

$$U^{j+1} = AU^j,$$

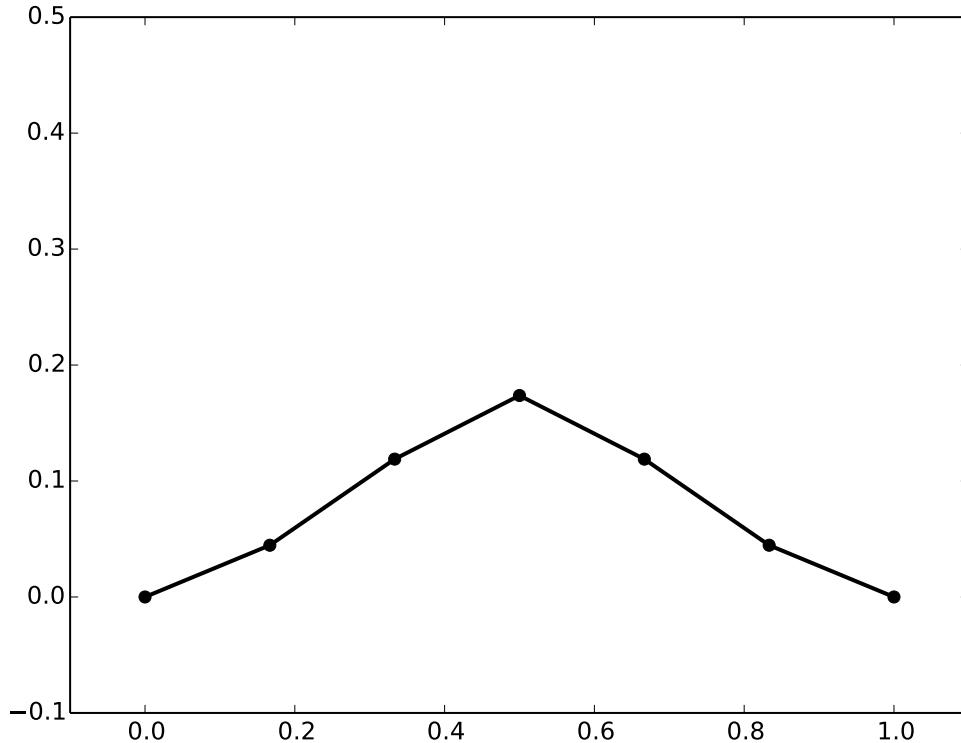


Figure 10.2: The graph of  $U^4$ , the approximation to the solution  $u(x, t = 0.4)$  for Problem 1.

where  $A$  is the tridiagonal matrix given by

$$A = \begin{bmatrix} 1 & 0 & & & \\ \lambda & 1 - 2\lambda & \lambda & & \\ & \ddots & \ddots & \ddots & \\ & & \lambda & 1 - 2\lambda & \lambda \\ & & & 0 & 1 \end{bmatrix},$$

$\lambda = \nu k / h^2$ , and  $U^j$  represents the approximation at time  $t_j$ . We can initialize this method using the initial condition given in our problem, which tells us that  $U_i^0 = f(x_i)$ .

Note

Note that the matrix representing the finite difference scheme is very sparse, which is typical of finite difference schemes. While representing finite difference schemes with matrices can be an effective method, especially for implicit schemes, it is very important that a sparse matrix format is used. Otherwise, performance will be dramatically negatively impacted. In Python, since looping is slow, the best alternative is to vectorize the difference scheme. This approach can in fact be even better than using matrices for explicit schemes, such as the one we are using here, as it avoids needing to store the matrix in memory.

To account for our constant boundary conditions using this differencing scheme, simply set the boundary points to the appropriate values in the initial conditions, then avoid modifying them as you update for each time step. Note that the first and last rows of the matrix representation of the differencing scheme are the same as the first and last rows of the identity matrix. This has the effect of keeping the boundary points the same as in the previous step, and thus the same as in the initial condition.

**Problem 1.** Consider the initial/boundary value problem

$$\begin{aligned} u_t &= 0.05u_{xx}, \quad x \in [0, 1], \quad t \in [0, 1] \\ u(0, t) &= 0, \quad u(1, t) = 0, \\ u(x, 0) &= 2 \max\{0.2 - |x - 0.5|, 0\}. \end{aligned} \tag{10.2}$$

Approximate the solution  $u(x, t)$  by taking 6 subintervals in the  $x$  dimension and 10 subintervals in time. Plot the solution at the times  $t = 0$ ,  $t = 0.4$ , and  $t = 1$ . The graphs for  $U^0$  and  $U^4$  are given in Figures 10.1 and 10.2.

**Problem 2.** Solve the initial/boundary value problem

$$\begin{aligned} u_t &= u_{xx}, \quad x \in [-12, 12], \quad t \in [0, 1], \\ u(-12, t) &= 0, \quad u(12, t) = 0, \\ u(x, 0) &= \max\{1 - x^2, 0\} \end{aligned} \tag{10.3}$$

using the first order explicit method 10.1. Use 140 subintervals in the  $x$  dimension and 70 subintervals in time. The initial and final states are shown in Figure 10.3. Animate your results.

Explicit methods usually have a stability condition, called a CFL condition (for Courant-Friedrichs-Lowy). For method 10.1 the CFL condition that must be satisfied is that

$$\lambda = \frac{\nu k}{h^2} \leq \frac{1}{2}.$$

Repeat your computations using 140 subintervals in the  $x$  dimension and 66 subintervals in time. Animate the results. For these values, the CFL condition is broken; you should be able to clearly see the result of this instability in the approximation  $U^{66}$ .

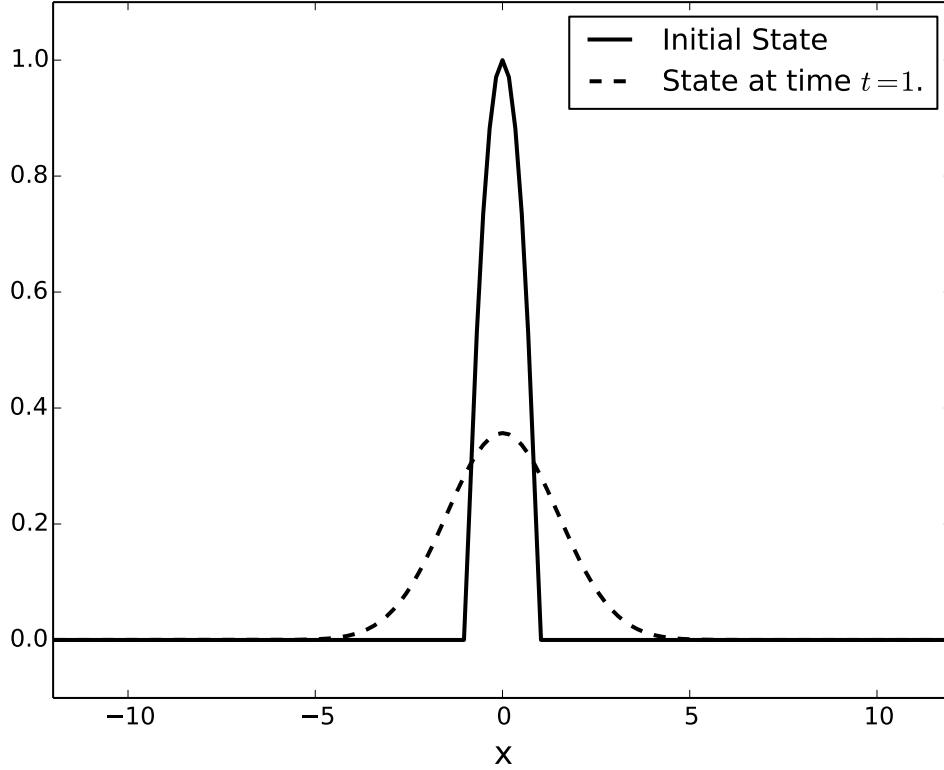


Figure 10.3: The initial and final states for equation Problem 2.

Implicit methods often have better stability properties than explicit methods. The Crank-Nicolson method, for example, is unconditionally stable and has order  $\mathcal{O}(h^2 + k^2)$ . To derive the Crank-Nicolson method, we use the following approximations:

$$\begin{aligned} u_t(x_i, t_{j+1/2}) &= \frac{u(x_i, t_{j+1}) - u(x_i, t_j)}{k} + \mathcal{O}(k^2), \\ u_{xx}(x_i, t_{j+1/2}) &= \frac{u_{xx}(x_i, t_{j+1}) + u_{xx}(x_i, t_j)}{2} + \mathcal{O}(k^2). \end{aligned}$$

The first equation is a finite difference approximation for  $u_t$ , and the second is a midpoint approximation applied to  $u_{xx}$ . These approximations give the relation

$$\begin{aligned} \frac{U_i^{j+1} - U_i^j}{k} &= \frac{1}{2} \left( \frac{U_{i+1}^j - 2U_i^j + U_{i-1}^j}{h^2} + \frac{U_{i+1}^{j+1} - 2U_i^{j+1} + U_{i-1}^{j+1}}{h^2} \right), \\ U_i^{j+1} &= U_i^j + \frac{k}{2h^2} \left( U_{i+1}^j - 2U_i^j + U_{i-1}^j + U_{i+1}^{j+1} - 2U_i^{j+1} + U_{i-1}^{j+1} \right). \end{aligned} \tag{10.4}$$

This method can be written in matrix form as

$$BU^{j+1} = AU^j,$$

where  $A$  and  $B$  are tridiagonal matrices given by

$$B = \begin{bmatrix} 1 & 0 & & & \\ -\lambda & 1+2\lambda & -\lambda & & \\ & \ddots & \ddots & \ddots & \\ & & -\lambda & 1+2\lambda & -\lambda \\ & & & 0 & 1 \end{bmatrix},$$

$$A = \begin{bmatrix} 1 & 0 & & & \\ \lambda & 1-2\lambda & \lambda & & \\ & \ddots & \ddots & \ddots & \\ & & \lambda & 1-2\lambda & \lambda \\ & & & 0 & 1 \end{bmatrix},$$

where  $\lambda = \nu k / (2h^2)$ , and  $U^j$  represents the approximation at time  $t_j$ . Note that here we have defined  $\lambda$  differently than we did before!

How do we know if a numerical approximation is reasonable? One way to determine this is to compute solutions for various step sizes  $h$  and see if the solutions are converging to something, which we hope to be the true solution. To be more specific, suppose our finite difference method is  $\mathcal{O}(h^p)$  accurate. This means that the error  $E(h) \approx Ch^p$  for some constant  $C$  as  $h \rightarrow 0$  (that is, for  $h > 0$  small enough).

So, we will compute the approximation  $y_k$  for each stepsize  $h_k$ ,  $h_1 > h_2 > \dots > h_m$ . We will think of  $y_m$  as the true solution. Then the error of the approximation for stepsize  $h_k$ ,  $k < m$ , is

$$E(h_k) = \max(|y_k - y_m|) \approx Ch_k^p,$$

$$\log(E(h_k)) = \log(C) + p \log(h_k).$$

Thus on a log-log plot of  $E(h)$  vs.  $h$ , these values should be on a straight line with slope  $p$  when  $h$  is small enough to start getting convergence.

**Problem 3.** Using the Crank Nicolson method, numerically approximate the solution  $u(x, t)$  of the problem

$$\begin{aligned} u_t &= u_{xx}, \quad x \in [-12, 12], \quad t \in [0, 1], \\ u(-12, t) &= 0, \quad u(12, t) = 0, \\ u(x, 0) &= \max\{1 - x^2, 0\}. \end{aligned} \tag{10.5}$$

Note that this is an implicit linear scheme; hence, the most efficient way to find  $U^{j+1}$  is to create the matrix  $B$  as a sparse matrix and use `scipy.sparse.linalg.spsolve`.

Demonstrate that the numerical approximation at  $t = 1$  converges. Do this by computing  $U$  at  $t = 1$  using 20, 40, 80, 160, 320, and 640 steps. Use the same number of steps in both time and space. Reproduce the loglog plot shown in Figure 10.4. The slope of the line there shows the order of convergence.

To measure the error, use the solution with the smallest  $h$  (largest number of intervals) as if it were the exact solution, then sample each solution only at the x-values that are represented in the solution with the largest  $h$  (smallest number of intervals). Use the  $\infty$ -norm on the arrays of values at those points to measure the error.

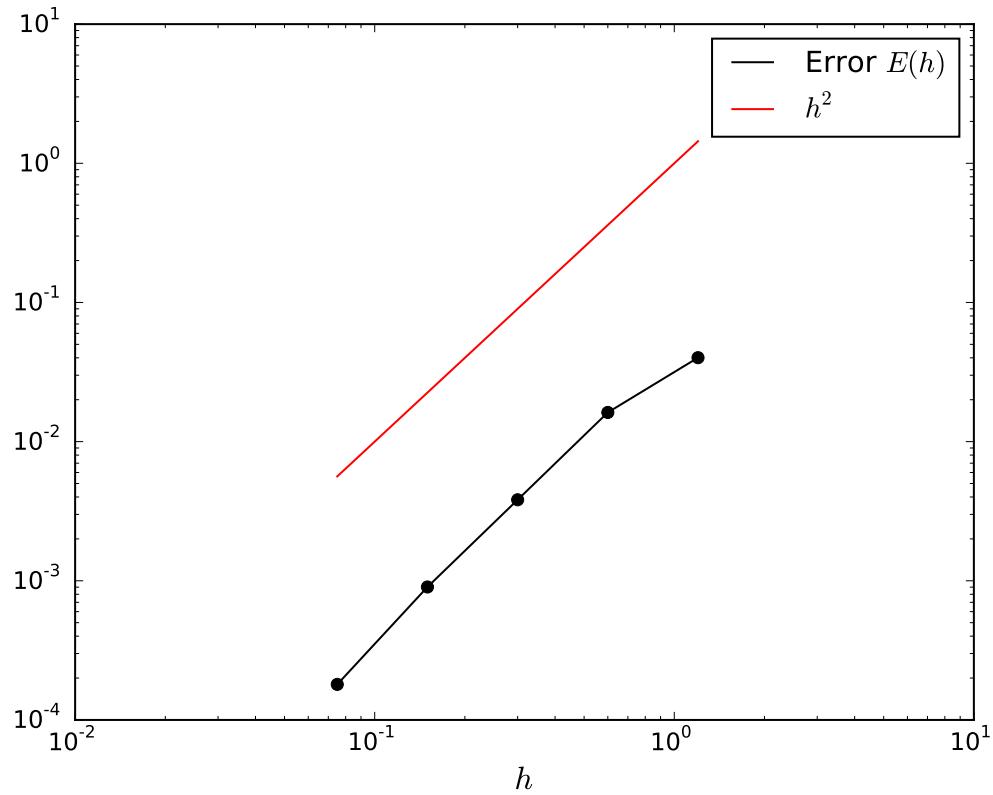


Figure 10.4:  $E(h)$  represents the (approximate) maximum error in the numerical solution  $U$  to Problem 3 at time  $t = 1$ , using a stepsize of  $h$ .

Notice that, since the Crank-Nicolson method is unconditionally stable, there is no CFL condition, and we can safely use the same number of intervals in time and space.

# 11

## Anisotropic Diffusion

**Lab Objective:** Demonstrate the use of finite difference schemes in image analysis.

A common task in image processing is to remove extra static from an image. This is most easily done by simply blurring the image, which can be accomplished by treating the image as a rectangular domain and applying the diffusion (heat) equation:

$$u_t = c\Delta u$$

where  $c$  is some diffusion constant and  $\Delta$  is the Laplace operator. Unfortunately, this also blurs the boundary lines between distinct elements of the image.

A more general form of the diffusion equation in two dimensions is:

$$u_t = \nabla \cdot (c(x, y, t)\nabla u)$$

where  $c$  is a function representing the diffusion coefficient at each given point and time. In this case,  $\nabla \cdot$  is the divergence operator and  $\nabla$  is the gradient.

To blur a picture uniformly, choose  $c$  to be a constant function. Since  $c$  controls how much diffusion is allowed at each point, it can be modified so that diffusion is minimized across edges in the image. In this way we attempt to limit diffusion near the boundaries between different features of the image, and allow smaller details of the image (such as static) to blur away. This method for image denoising is especially useful for denoising low quality images, and was first introduced by Pietro Perona and Jitendra Malik in 1987. It is known as Anisotropic Diffusion or Perona-Malik Diffusion.

### A Finite Difference Scheme

Suppose we have some estimate  $E$  of the rate of change at a given point in an image.  $E$  will be largest at the boundaries in the image. We will then let  $c(x, y, t) = g(E(x, y, t))$  where  $g$  is some function such that  $g(0) = 1$  and  $\lim_{x \rightarrow \infty} g(x) = 0$ . Thus  $c$  will be small where  $E$  is large, so that little diffusion occurs near the boundaries of different portions of the image.

We will model this system using a finite differencing scheme with an array of values at a 2D grid of points, and iterate through time. Let  $U_{l,m}^n$  be the discretized approximation of the function  $u$ ,  $n$  be the index in time,  $l$  be the index along the  $x$ -axis, and  $m$  be the index along the  $y$ -axis.

The Laplace operator can be approximated with the finite difference scheme

$$\Delta u = u_{xx} + u_{yy} \approx \frac{U_{l-1,m}^n - 2U_{l,m}^n + U_{l+1,m}^n}{(\Delta x)^2} + \frac{U_{l,m-1}^n - 2U_{l,m}^n + U_{l,m+1}^n}{(\Delta y)^2}.$$

A good metric to use with images is to let the distance between each pixel be equal to one, so  $\Delta x = \Delta y = 1$ . Rearranging terms, we obtain

$$\Delta u \approx (U_{l-1,m}^n - U_{l,m}^n) + (U_{l+1,m}^n - U_{l,m}^n) + (U_{l,m-1}^n - U_{l,m}^n) + (U_{l,m+1}^n - U_{l,m}^n).$$

Again, since we are working with images and not some time based problem, we can without loss of generality let  $\Delta t = 1$ , so we obtain the finite difference scheme

$$U_{l,m}^{n+1} = U_{l,m}^n + (U_{l-1,m}^n - U_{l,m}^n) + (U_{l+1,m}^n - U_{l,m}^n) + (U_{l,m-1}^n - U_{l,m}^n) + (U_{l,m+1}^n - U_{l,m}^n).$$

We will now limit the diffusion near the edges of objects by making the modification

$$\begin{aligned} U_{l,m}^{n+1} = U_{l,m}^n &+ \lambda \left( g(|U_{l-1,m}^n - U_{l,m}^n|)(U_{l-1,m}^n - U_{l,m}^n) \right. \\ &+ g(|U_{l+1,m}^n - U_{l,m}^n|)(U_{l+1,m}^n - U_{l,m}^n) \\ &+ g(|U_{l,m-1}^n - U_{l,m}^n|)(U_{l,m-1}^n - U_{l,m}^n) \\ &\left. + g(|U_{l,m+1}^n - U_{l,m}^n|)(U_{l,m+1}^n - U_{l,m}^n) \right), \end{aligned} \quad (11.1)$$

where  $\lambda \leq \frac{1}{4}$  is the stability condition.

In this difference scheme, each term is affected most by nearby terms that are most similar to it, so less diffusion will happen anywhere there is a sharp difference between pixels. This scheme also has the useful property that it does not increase or decrease the total brightness of the image. Intuitively, this is because the effect of each point on its neighbors is exactly the opposite effect its neighbors have on it.

Two commonly used functions for  $g$  are  $g(x) = e^{-(\frac{x}{\sigma})^2}$  and  $g(x) = \frac{1}{1+(\frac{x}{\sigma})^2}$ . The parameter  $\sigma$  allows us to control how much diffusion decreases across boundaries, with larger  $\sigma$  values allowing more diffusion. Note that  $g(0) = 1$  and  $\lim_{x \rightarrow \infty} g(x) = 0$  for both functions. In this lab we use  $g(x) = e^{-(\frac{x}{\sigma})^2}$ .

It is worth noting that this particular difference scheme is not an accurate finite difference scheme for the version of the diffusion equation we discussed before, but it does accomplish the same thing in the same way. As it turns out, this particular scheme is the solution to a slightly different diffusion PDE, but can still be used the same way.

For this lab's examples we read in the image using the `imageio.imread` function, and normalized it so that the colors are represented as floating point values between 0 and 1. An image can be converted to black and white when it is read by including the argument `as_gray=True`.

```
from matplotlib import cm, pyplot as plt
from imageio import imread

# To read in an image, convert it to grayscale, and rescale it.
picture = imread('balloon.jpg', as_gray=True) * 1./255

# To display the picture as grayscale
plt.imshow(picture, cmap=cm.gray)
plt.show()
```

## Simplifying Calculations

You will notice that the algorithm given in 11.1 does not describe what to do for the edges and corners of  $U^{n+1}$ . In these cases we will simply eliminate the undefined terms in the algorithm. For example, the top edge equation becomes

$$\begin{aligned} U_{l,m}^{n+1} = & U_{l,m}^n + \lambda(g(|U_{l+1,m}^n - U_{l,m}^n|)(U_{l+1,m}^n - U_{l,m}^n) \\ & + g(|U_{l,m+1}^n - U_{l,m}^n|)(U_{l,m+1}^n - U_{l,m}^n)) \\ & + g(|U_{l,m-1}^n - U_{l,m}^n|)(U_{l,m-1}^n - U_{l,m}^n)), \end{aligned}$$

and top left corner equation becomes

$$\begin{aligned} U_{l,m}^{n+1} = & U_{l,m}^n + \lambda(g(|U_{l+1,m}^n - U_{l,m}^n|)(U_{l+1,m}^n - U_{l,m}^n) \\ & + g(|U_{l,m+1}^n - U_{l,m}^n|)(U_{l,m+1}^n - U_{l,m}^n)). \end{aligned}$$

Essentially we are only using the terms of the difference scheme that are actually defined.

To help facilitate this we can create a larger "padded" matrix that will make these calculations easy to do. This padded matrix will have an extra row on the top and bottom, and an extra column on either side of the original matrix. These extra rows and columns will duplicate the outer edge of the original matrix.

So if our original array  $X$  has shape  $m,n$ , then our padded array  $Y$  has shape  $m+2,n+2$ . The top edge of  $Y$  will be defined so that  $Y[0,1:-1] == X[0,:]$  is true, and the rest of the edges of  $Y$  follow the same pattern.

Notice that this allows us to simply implement the algorithm found in 11.1 without having to make special cases for the edges and corners, since those previously undefined terms become zero when using the padded matrix.

**Problem 1.** Complete the following function, by implementing the anisotropic diffusion algorithm found in 11.1 for black and white images. Use the padded array technique found in the Simplifying Calculations section.

In your function, use

$$g(x) = e^{-(\frac{x}{\sigma})^2}$$

```
def anisdiff_bw(U, N, lambda_, g):
    """ Run the Anisotropic Diffusion differencing scheme
    on the array U of grayscale values for an image.
    Perform N iterations, use the function g
    to limit diffusion across boundaries in the image.
    Operate on U inplace to optimize performance. """
    pass
```

Run the function on `balloon.jpg`. Show the original image and the diffused image for  $\sigma = .1$ ,  $\lambda = .25$ ,  $N = 5, 20, 100$ .



original image

5 iterations with  $\sigma = .1$  and  $\lambda = .25$ 

20 iterations



100 iterations

## Color Schemes

Colored images can be processed in a similar manner. Instead of being represented as a two-dimensional array, colored images are represented as three dimensional arrays. The third dimension is used to store the intensities of each of the standard 3 colors. This diffusion process can be carried out in the exact same way, on each of the arrays of intensities for each color, but instead of detecting edges just in one color, we need to detect edges in any color, so instead of using something of the form  $g(|U_{l+1,m}^n - U_{l,m}^n|)$  as before, we will now use something of the form  $g(||U_{l+1,m}^n - U_{l,m}^n||)$ , where  $U_{l+1,m}^n$  and  $U_{l,m}^n$  are vectors now instead of scalars. The difference scheme can be treated as an equation on vectors in 3-space and now reads:

$$\begin{aligned} U_{l,m}^{n+1} = & U_{l,m}^n + \lambda(g(||U_{l-1,m}^n - U_{l,m}^n||)(U_{l-1,m}^n - U_{l,m}^n) \\ & + g(||U_{l+1,m}^n - U_{l,m}^n||)(U_{l+1,m}^n - U_{l,m}^n) \\ & + g(||U_{l,m-1}^n - U_{l,m}^n||)(U_{l,m-1}^n - U_{l,m}^n) \\ & + g(||U_{l,m+1}^n - U_{l,m}^n||)(U_{l,m+1}^n - U_{l,m}^n)) \end{aligned}$$

When implementing this scheme for colored images, use the 2-norm on 3-space, i.e  $\|x\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$  where  $x_1$ ,  $x_2$ , and  $x_3$  are the different coordinates of  $x$ .

**Problem 2.** Complete the following function to process a colored image. You may modify your code from the previous problem. Measure the difference between pixels using the 2-norm. Use the corresponding vector versions of the boundary conditions given in Problem 1.

```
def anisdiff_color(U, N, lambda_, g):
    """ Run the Anisotropic Diffusion differencing scheme
    on the array U of color values for an image.
    Perform N iterations, use the function g = e^{-x^2/sigma^2}
    to limit diffusion across boundaries in the image.
    Operate on U inplace to optimize performance. """
    pass
```

Run the function on `balloons_color.jpg`. Show the original image and the diffused image for  $\sigma = .1$ ,  $\lambda = .25$ ,  $N = 5, 20, 100$ .

Hint: If you have an  $m \times n \times 3$  matrix representing the RGB differences of each pixel, then to find a matrix representing the norm of the differences, you can use the following code. This code squares each value and sums along the last axis, and takes the square root. In order to keep the dimension size of the matrix and aid in broadcasting, you must use `keepdims=True`.

```
# x is mxnx3 matrix of pixel color values
norm = np.sqrt(np.sum(x**2, axis=2, keepdims=True))
```

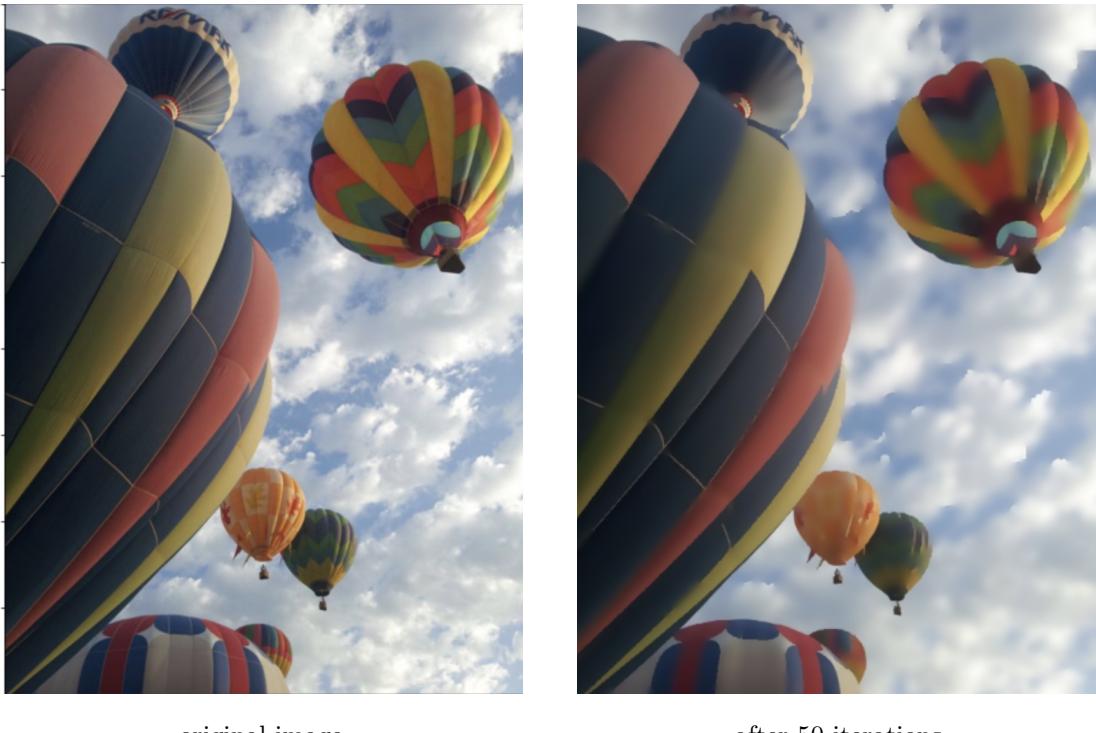


Figure 11.1: Smearing of similar colors when using an anisotropic diffusion filter.

## Noisy Images

**Problem 3.** Use the following code to add noise to your grayscale image.

```
from numpy.random import randint

image = imread('balloon.jpg', as_gray=True)
x, y = image.shape
for i in range(x*y//100):
    image[randint(x),randint(y)] = 127 + randint(127)
```

Run `anisdiff_bw()` on the noisy image with  $\sigma = .1$ ,  $\lambda = .25$ ,  $N = 20$ . Display the original image and the noisy image. Explain why anisotropic diffusion does not smooth out the noise.

Hint: Don't forget to rescale.

## Minimum Bias

This sort of anisotropic diffusion can be very effective, but, depending on the image, it may also smear out edges that do not have large differences between them. An example of this limitation can be seen in Figure 11.1

As we can see, after 100 iterations, some of the boundaries between similar shades of grey have smeared unevenly. You may still have to look closely to see it. This can be counteracted somewhat by further decreasing the  $\sigma$  value, but if we have random noise throughout the image, this will not remove it. If we have random static in the image, we can remove this using a modified version of the filter. Instead of measuring the rate of change in the picture in each direction, we change each point according to whether or not any of its adjacent points have roughly the same value it has. This is called a minimum-biased filter. This sort of trick is especially good for removing isolated pixels that are different from those around them. A very simple way to do this is by taking the average of the two smallest differences between each pixel and its eight neighbors and using that in place of  $g$  in the difference scheme above. Along the boundaries, we do not have 8 neighbors for each pixel, but we can get by just using the pixels we have and eliminating the other terms in the difference scheme, just as we did before. This will make it so that points that neighbor points of similar value will not be changed, while points that do not match their surroundings will be faded to become more like the points surrounding them. This does not have the same symmetrical diffusion as the other scheme, i.e. if one pixel changes, it does not necessarily change its neighboring pixels by the same amount. As long as you leave  $\lambda \leq \frac{1}{4}$  and you have scaled the pixels to have floating point values between 0 and 1, the scheme will still remain within its minimum and maximum bounds, since the tendency is always to move points closer to the values of their neighbors. To demonstrate the action of such a filter, we make changes to random pixels in the color version of the same photo and use both filters to remove the noise we have added. Below, we include an example where we have added noise to the color version of that same picture, then used a minimum-biased filter to diminish the noise and the original filter to smooth what remains.

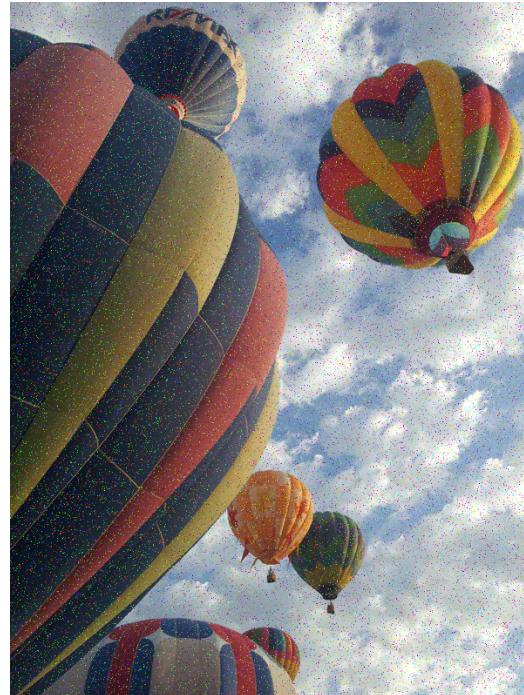
**Problem 4.** Implement the minimum-biased finite difference scheme described above. Add noise to `balloons_color.jpg` using the provided code below, and clean it using your implementation. Show the original image, the noised image, and the cleaned image.

```
image = imread('balloons_color.jpg')
x,y,z = image.shape
for dim in range(z):
    for i in range(x*y//100):
        # Assign a random value to a random place
        image[randint(x),randint(y),dim] = 127 + randint(127)
```

Hint: Don't forget to rescale.



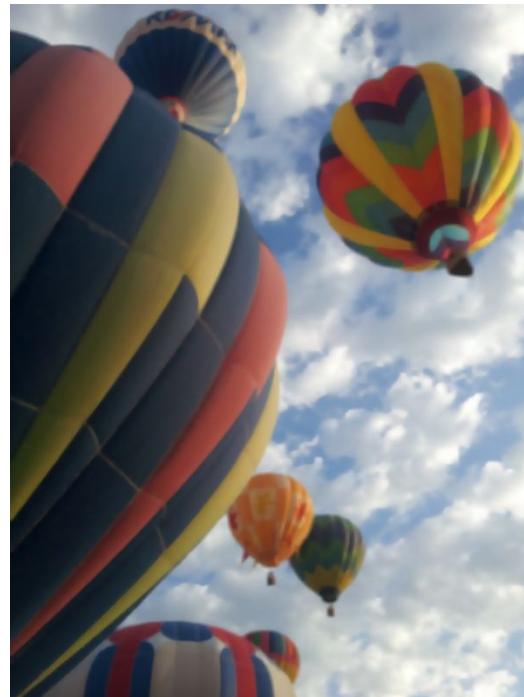
original image



randomly changed 100000 color values



300 iterations of a min-biased scheme

after 8 additional iterations of the first filter  
with  $\lambda = .25$  and  $\sigma = .04$ .

# 12

## The Finite Element method

**Lab Objective:** The finite element method is commonly used for numerically solving partial differential equations. We introduce the finite element method via a simple BVP describing the steady state distribution of heat in a pipe as fluid flows through.

### Advection-Diffusion of Heat in a Fluid

We wish to study the distribution of heat in a fluid that is moving at some constant speed  $a$ . Let  $y$  denote the temperature of the fluid at any given location and time. The equation modeling this situation can be obtained from the differential form of the conservation law, where the flux is the sum of a diffusive term  $-\varepsilon y_x$  and an advection (or transport) term  $ay$ :

$$J = ay - \varepsilon y_x$$

The one-dimension conservation law states that  $y$  must then obey the partial differential equation  $y_t + J_x = f(x)$ , where  $f$  represents heat sources in the system. Since  $J_x = ay_x - \varepsilon y_{xx}$ , we obtain the advection-diffusion equation

$$y_t + ay_x = \varepsilon y_{xx} + f(x).$$

As time progresses, we expect the temperature of the fluid in the pipe to reach a steady state distribution, with  $y_t = 0$ . Once this steady state has been reached, the heat distribution  $y$  then satisfies the ODE

$$\varepsilon y'' - ay' = -f(x).$$

We consider the scenario of a fluid flowing through a pipe from  $x = 0$  to  $x = 1$  with speed  $a = 1$ , and as it travels it is warmed at a constant rate  $f(x) = 1$ . Note that since this a second-order ODE, we need two boundary conditions. Suppose that the fluid is already at a known temperature  $y = 2$  as it enters the pipe. This imposes the boundary condition  $y(0) = 2$ . Suppose further that a device is installed on the end of the pipe that nearly instantaneously brings the heat of the water up to  $y = 4$ . Physically, we expect this extra heat that is introduced at  $x = 1$  to diffuse backward through the water in the pipe and thus influence the steady-state temperature. Putting this together leads to a well defined BVP:

$$\begin{aligned} \varepsilon y'' - y' &= -1, & 0 < x < 1, \\ y(0) &= 2, & y(1) = 4. \end{aligned} \tag{12.1}$$

The analytic solution for  $\varepsilon = 0.1$  is shown in Figure 12.1.

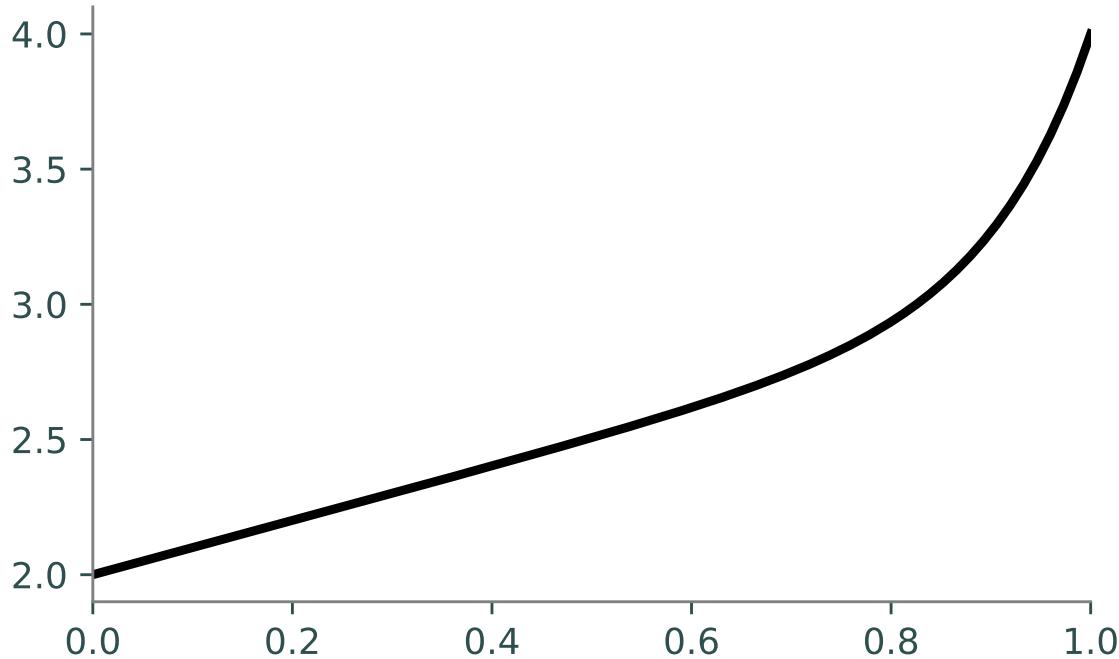


Figure 12.1: The analytic solution of (12.1) for  $\varepsilon = 0.1$ .

## The Weak Formulation

Stepping back momentarily, consider the equation

$$\begin{aligned} \varepsilon y'' - y' &= -f, \quad 0 < x < 1, \\ y(0) &= \alpha, \quad y(1) = \beta. \end{aligned} \tag{12.2}$$

To approximate the solution  $y$  using the finite element method, we reframe the problem into one involving integrals, known as its weak formulation.

Let  $w$  be a smooth function on  $[0, 1]$  satisfying  $w(0) = w(1) = 0$ . Multiplying (12.2) by  $w$  and integrating over  $[0, 1]$  yields

$$\begin{aligned} \int_0^1 -fw \, dx &= \int_0^1 (\varepsilon y'' w - y' w) \, dx, \\ &= \int_0^1 (-\varepsilon y' w' - y' w) \, dx, \end{aligned}$$

where the second equality follows by integration by parts. For notational convenience, define the functionals  $a$  and  $l$  by

$$\begin{aligned} a(y, w) &= \int_0^1 -\varepsilon y' w' - y' w \, dx, \\ l(w) &= \int_0^1 -fw \, dx. \end{aligned}$$

Then, any solution to (12.2) will also satisfy

$$a(y, w) = l(w) \tag{12.3}$$

This equation is the weak formulation of (12.2). Note that any solution to the original ODE is also a solution to the weak formulation. However, solutions to the weak formulation need not be solutions to the original ODE, as they may not even be differentiable everywhere. While this may seem like an undesirable property, it allows us to use a wider variety of functions to approximate the true solution.

Now, we choose some appropriate vector space  $V$  of functions, and consider the problem of finding a function  $y \in V$  that satisfies the weak formulation (12.3) for all  $w \in V_0 = \{w \in V | w(0) = w(1) = 0\}$ . The finite element method consists of choosing  $V$  to be some set of piecewise polynomial functions. In this lab, we will consider the case of using piecewise linear functions.

## The Finite Element Method

Let  $P_n$  be some partition of  $[0, 1]$ ,  $0 = x_0 < x_1 < \dots < x_n = 1$ , and let  $V_n$  be the set of continuous linear piecewise functions  $v$  on  $[0, 1]$  such that  $v$  is linear on each subinterval  $[x_j, x_{j+1}]$ . These subintervals are the finite elements for which this method is named. Note that  $V_n$  has dimension  $n + 1$ , since each of the continuous piecewise linear functions in  $V$  are uniquely determined by their values at the  $n + 1$  points  $x_0, x_1, \dots, x_n$ . Let  $V_{n,0}$  be the subspace of  $V_n$  of dimension  $n - 1$  whose elements are zero at the endpoints of  $[0, 1]$ .

Let the  $\phi_i$  be the hat functions

$$\phi_i(x) = \begin{cases} (x - x_{i-1})/h_i & \text{if } x \in [x_{i-1}, x_i] \\ (x_{i+1} - x)/h_{i+1} & \text{if } x \in [x_i, x_{i+1}] \\ 0 & \text{otherwise} \end{cases}$$

where  $h_i = x_i - x_{i-1}$ ; see Figures 12.2 and 12.3. These hat functions form a basis for  $V_n$ . Note that the points  $x_0, \dots, x_n$  need not be evenly spaced, and the  $h_i$  do not need to be equal. This is in fact one of the major strengths of this approach, as it allows adapting the points in the partition to the problem, which can reduce the error in the approximation. When applied to PDEs, it also is a simple way to handle unusually-shaped domains.

We now can write our approximate solution for  $y$  and the arbitrary function  $w$  as a linear combination of these basis elements, which will enable us to solve the system numerically. In particular, we can write  $\hat{y}(x) = \sum_{i=0}^n k_i \phi_i(x)$ , where the  $k_i$  are to be determined.

To make things more concrete, consider the case of  $n = 5$  with the partition  $P_5 = \{x_0, x_1, \dots, x_5\}$ . We look for an approximation  $\hat{y} = \sum_{i=0}^5 k_i \phi_i \in V_5$  of the true solution  $y$ ; to do this, we must determine appropriate values for the constants  $k_i$ . We impose the condition on  $\hat{y}$  that

$$a(\hat{y}, w) = l(w)$$

for all  $w \in V_{5,0}$ . This can be written equivalently as

$$a\left(\sum_{i=0}^5 k_i \phi_i, \phi_j\right) = l(\phi_j) \quad \text{for } j = 1, 2, 3, 4,$$

since  $a$  and  $l$  are linear in  $w$  and  $\phi_1, \phi_2, \phi_3, \phi_4$  form a basis for  $V_{5,0}$ . Since  $a$  is also linear in  $y$ , we further obtain

$$\sum_{i=0}^5 k_i a(\phi_i, \phi_j) = l(\phi_j) \quad \text{for } j = 1, 2, 3, 4.$$

To satisfy the boundary conditions, we necessarily have that  $k_0 = \alpha$ ,  $k_5 = \beta$ . These equations can be written together in matrix form as

$$AK = \Phi, \tag{12.4}$$

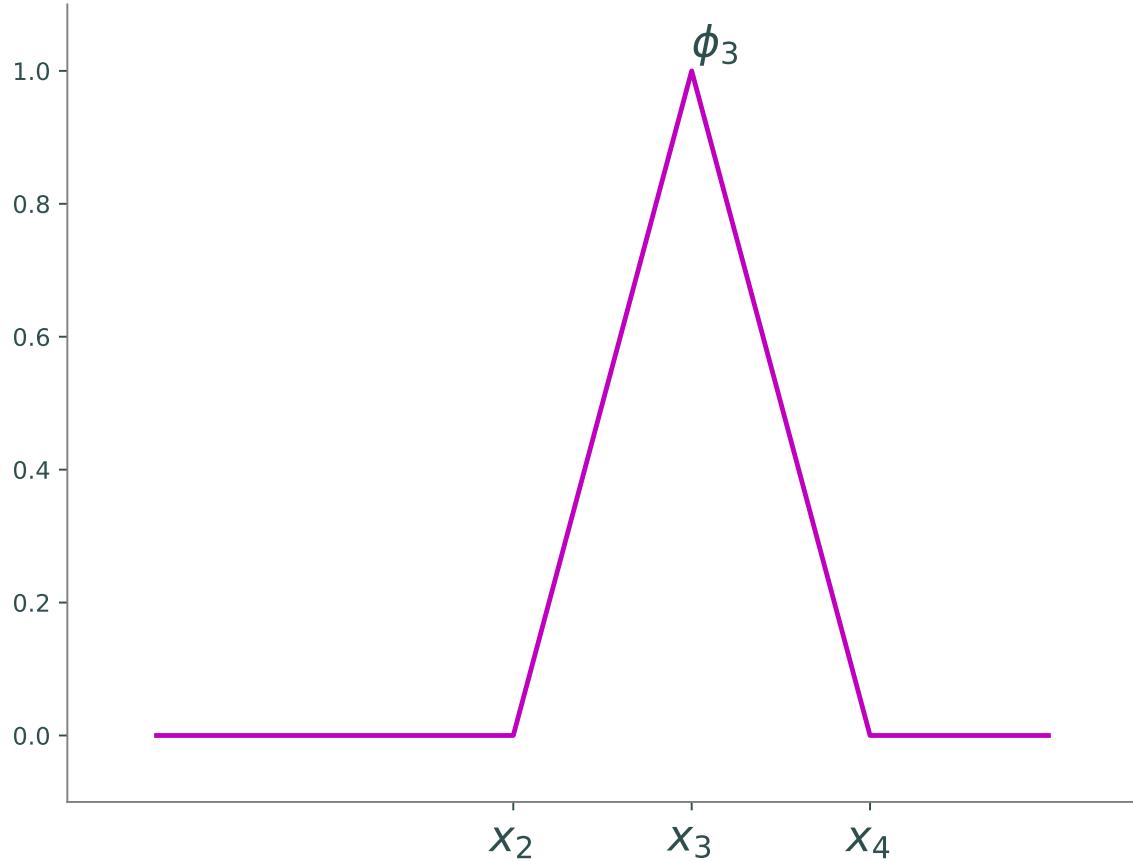


Figure 12.2: The basis function  $\phi_3$ , when the  $x_i$  are evenly spaced.

where

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ a(\phi_0, \phi_1) & a(\phi_1, \phi_1) & a(\phi_2, \phi_1) & 0 & 0 & 0 \\ 0 & a(\phi_1, \phi_2) & a(\phi_2, \phi_2) & a(\phi_3, \phi_2) & 0 & 0 \\ 0 & 0 & a(\phi_2, \phi_3) & a(\phi_3, \phi_3) & a(\phi_4, \phi_3) & 0 \\ 0 & 0 & 0 & a(\phi_3, \phi_4) & a(\phi_4, \phi_4) & a(\phi_5, \phi_4) \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

and

$$K = \begin{bmatrix} k_0 \\ k_1 \\ k_2 \\ k_3 \\ k_4 \\ k_5 \end{bmatrix}, \quad \Phi = \begin{bmatrix} \alpha \\ l(\phi_1) \\ l(\phi_2) \\ l(\phi_3) \\ l(\phi_4) \\ \beta \end{bmatrix}.$$

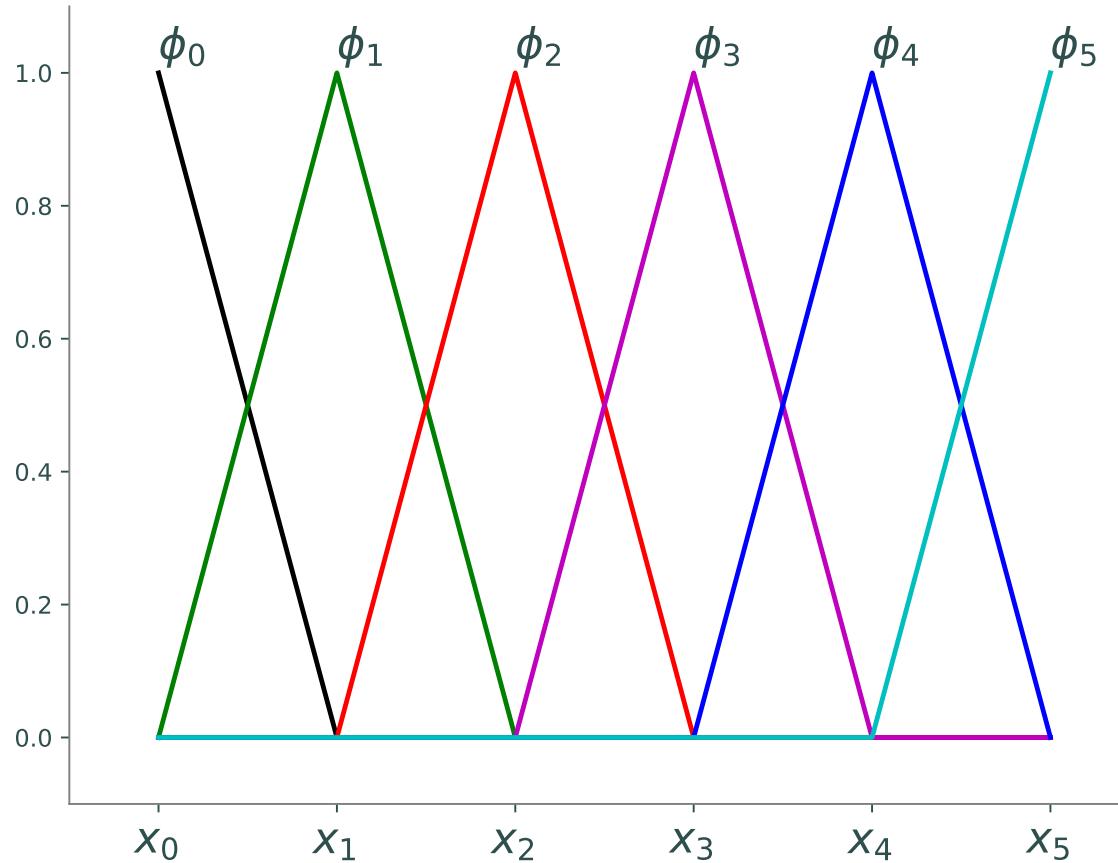


Figure 12.3: The six basis functions for  $V_5$ , when the  $x_i$  are evenly spaced.

Note that since  $a(\phi_i, \phi_j) = 0$  for most values of  $i, j$  (in particular, when the hat functions do not have overlapping domains), the finite element method results in a sparse linear system. To compute the coefficients of (12.4) we begin by evaluating some integrals. Since

$$\phi_i(x) = \begin{cases} (x - x_{i-1})/h_i & \text{if } x \in [x_{i-1}, x_i] \\ (x_{i+1} - x)/h_{i+1} & \text{if } x \in [x_i, x_{i+1}] \\ 0 & \text{otherwise} \end{cases}$$

$$\phi'_i(x) = \begin{cases} 1/h_i & \text{for } x_{i-1} < x < x_i, \\ -1/h_{i+1} & \text{for } x_i < x < x_{i+1}, \\ 0 & \text{otherwise,} \end{cases}$$

we obtain

$$\int_0^1 \phi'_i \phi'_j = \begin{cases} -1/h_{i+1} & \text{if } j = i + 1, \\ 1/h_i + 1/h_{i+1} & \text{if } j = i, \\ 0 & \text{otherwise,} \end{cases}$$

$$\int_0^1 \phi'_i \phi_j = \begin{cases} -1/2 & \text{if } j = i + 1, \\ 1/2 & \text{if } j = i - 1, \\ 0 & \text{otherwise,} \end{cases}$$

which can be put together to obtain (for  $f(x) = 1$ )

$$a(\phi_i, \phi_j) = \begin{cases} \varepsilon/h_{i+1} + 1/2 & \text{if } j = i + 1, \\ -\varepsilon/h_i - \varepsilon/h_{i+1} & \text{if } j = i, \\ \varepsilon/h_i - 1/2 & \text{if } j = i - 1, \\ 0 & \text{otherwise,} \end{cases}$$

$$l(\phi_j) = -\frac{1}{2}(h_j + h_{j+1}).$$

Equation (12.4) may now be solved using any standard linear solver. To handle the large number of elements required for Problem 3, you will want to use sparse matrices from `scipy.sparse`.

**Problem 1.** Use the finite element method to solve

$$\begin{aligned} \varepsilon y'' - y' &= -1, \\ y(0) = \alpha, \quad y(1) &= \beta, \end{aligned} \tag{12.5}$$

where  $\alpha = 2$ ,  $\beta = 4$ , and  $\varepsilon = 0.02$ . Use  $N = 100$  finite elements (101 grid points). Compare your solution with the analytic solution

$$y(x) = \alpha + x + (\beta - \alpha - 1) \frac{e^{x/\varepsilon} - 1}{e^{1/\varepsilon} - 1}.$$

Hint: One additional nice consequence of this setup is that the approximation  $\hat{y}$  is exactly the piecewise linear function that connects the points  $(x_i, k_i)$ . This means that the solution can be plotted very simply using `plt.plot(x, k)`, where `x` and `k` are arrays of the  $x_i$  and  $k_i$ .

**Problem 2.** One of the strengths of the finite element method is the ability to generate grids that better suit the problem. The solution of (12.5) changes most rapidly near  $x = 1$ . Compare the numerical solution when the grid points are unevenly spaced versus when the grid points are clustered in the area of greatest change; see Figure 12.4. Specifically, use the grid points defined by

```
even_grid = np.linspace(0,1,15)
clustered_grid = np.linspace(0,1,15)**(1./8)
```

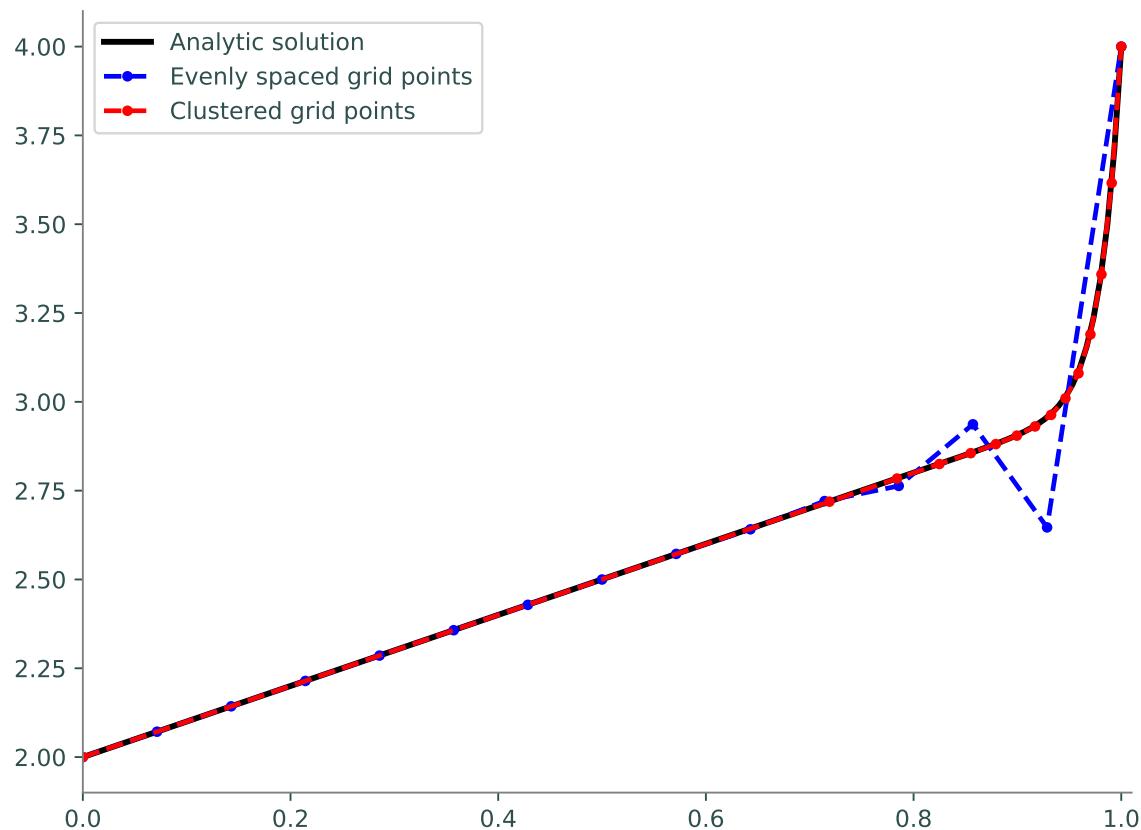


Figure 12.4: Two finite element approximations using 15 grid points, with different spacings.

**Problem 3.** Higher order methods promise faster convergence, but typically require more work to code. So why do we use them when a low order method will converge just as well, albeit with more grid points? The answer concerns the roundoff error associated with floating point arithmetic. Low order methods generally require more floating point operations, so roundoff error has a much greater effect.

The finite element method introduced here is a second order method, even though the approximate solution is piecewise linear. (To see this, note that if the grid points are evenly spaced, the matrix  $A$  in (12.4) is exactly the same as the matrix for the second order centered finite difference method.)

Solve (12.5) with the finite element method using  $N = 2^i$  evenly-spaced finite elements,  $i = 4, 5, \dots, 21$ . Remember to use sparse matrices, as this greatly reduces the memory and computation needed for the larger  $N$ . Compute the error as the maximum absolute value of the difference of the values of the approximate and true solutions at each of the  $x_i$ . Use a log-log plot to graph the error, and compare with Figure 12.5.

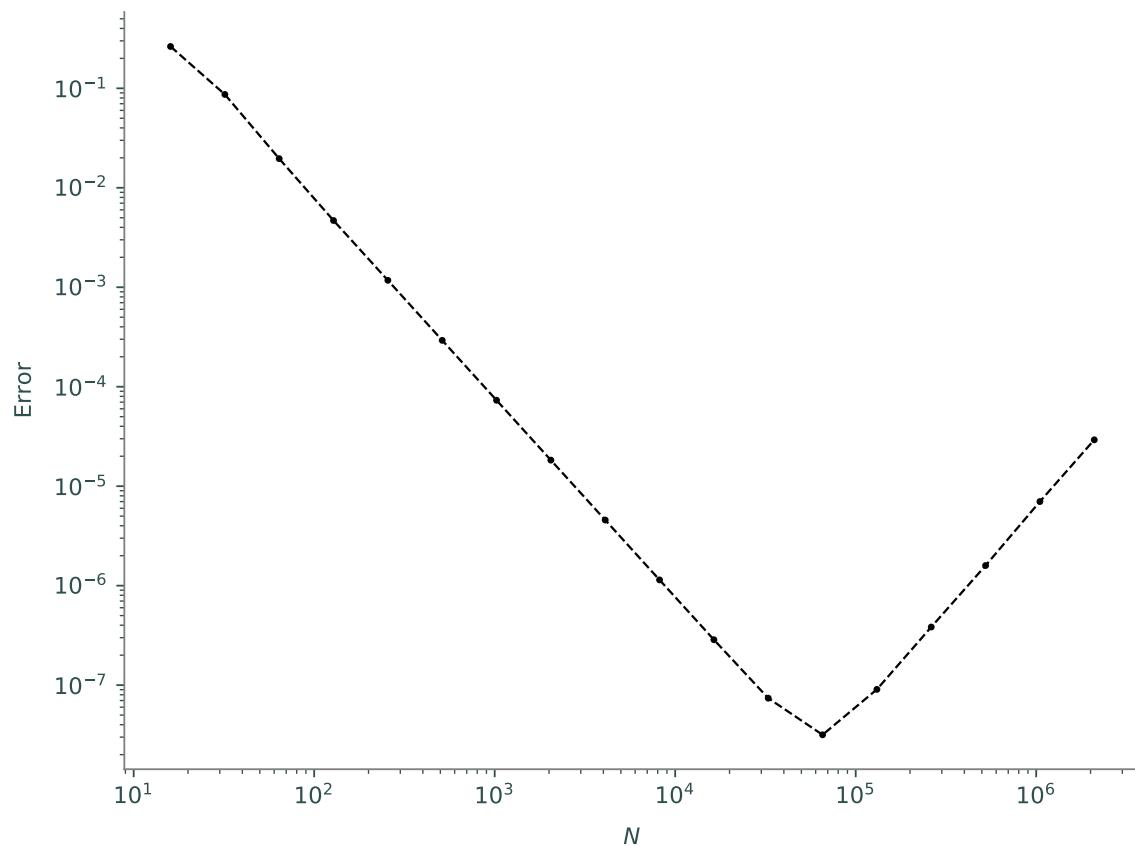


Figure 12.5: Error for the second order finite element method, as the number of subintervals  $N$  grows. Round-off error eventually overwhelms the approximation.

# 13 Poisson's equation

Suppose that we want to describe the distribution of heat throughout a region  $\Omega$ . Let  $h(x)$  represent the temperature on the boundary of  $\Omega$  ( $\partial\Omega$ ), and let  $g(x)$  represent the initial heat distribution at time  $t = 0$ . If we let  $f(x, t)$  represent any heat sources/sinks in  $\Omega$ , then the flow of heat can be described by the boundary value problem (BVP)

$$\begin{aligned} u_t &= \Delta u + f(x, t), \quad x \in \Omega, \quad t > 0, \\ u(x, t) &= h(x), \quad x \in \partial\Omega, \\ u(x, 0) &= g(x). \end{aligned} \tag{13.1}$$

When the source term  $f$  does not depend on time, there is often a steady-state heat distribution  $u_\infty$  that is approached as  $t \rightarrow \infty$ . This steady state  $u_\infty$  is a solution of the BVP

$$\begin{aligned} \Delta u + f(x) &= 0, \quad x \in \Omega, \\ u(x, t) &= h(x), \quad x \in \partial\Omega. \end{aligned} \tag{13.2}$$

This last partial differential equation,  $\Delta u = -f$ , is called Poisson's equation. This equation is satisfied by the steady-state solutions of many other evolutionary processes. Poisson's equation is often used in electrostatics, image processing, surface reconstruction, computational fluid dynamics, and other areas.

## Poisson's equation in two dimensions

Consider Poisson's equation together with Dirichlet boundary conditions on a rectangular domain  $R = [a, b] \times [c, d]$ :

$$\begin{aligned} u_{xx} + u_{yy} &= f, \quad x \text{ in } R \subset \mathbb{R}^2, \\ u &= g, \quad x \text{ on } \partial R. \end{aligned} \tag{13.3}$$

Let  $a = x_0, x_1, \dots, x_n = b$  and  $c = y_0, y_1, \dots, y_n = d$  be evenly spaced grids. Furthermore, suppose that  $b - a = d - c$ , so the rectangular domain is also square. Thus we have a single stepsize  $h$ , where  $h = x_{i+1} - x_i = y_{i+1} - y_i$

We look for an approximation  $U_{i,j}$  on the grid  $\{(x_i, y_j)\}_{i,j=0}^n$ .

Recall that

$$\begin{aligned}\Delta u &= u_{xx}(x, y) + u_{yy}(x, y) \\ &= \frac{u(x+h, y) - 2u(x, y) + u(x-h, y)}{h^2} \\ &\quad + \frac{u(x, y+h) - 2u(x, y) + u(x, y-h)}{h^2} + \mathcal{O}(h^2).\end{aligned}$$

We replace  $\Delta$  with the finite difference operator  $\Delta_h$ , defined by

$$\begin{aligned}\Delta_h U_{ij} &= \frac{U_{i+1,j} - 2U_{i,j} + U_{i-1,j}}{h^2} + \frac{U_{i,j+1} - 2U_{i,j} + U_{i,j-1}}{h^2}, \\ &= \frac{1}{h^2}(U_{i-1,j} + U_{i+1,j} + U_{i,j-1} + U_{i,j+1} - 4U_{i,j}).\end{aligned}$$

These equations are linear, so we can expect to write them in matrix form. However, since our unknown variables are doubly-indexed (for  $x_i$  and  $y_j$ ), we first need to rewrite them as a 1-dimensional array. We can do this by "stacking" the columns of the 2-dimensional array. Let the vector of unknowns  $U$  be:

$$U = \begin{bmatrix} U^1 \\ U^2 \\ \vdots \\ U^{n-1} \end{bmatrix} \text{ where } U^j = \begin{bmatrix} U_{1,j} \\ U_{2,j} \\ \vdots \\ U_{n-1,j} \end{bmatrix} \text{ for each } j, 1 \leq j \leq n-1.$$

Then the set of equations

$$\Delta_h U_{ij} = f_{ij}, \quad i, j = 1, \dots, n-1,$$

can be written in matrix form as

$$AU + p + q = f. \quad (13.4)$$

$A$  is a block tridiagonal matrix, given by

$$\frac{1}{h^2} \begin{bmatrix} T & I & & & \\ I & T & I & & \\ & \ddots & \ddots & \ddots & \\ & & I & T & I \\ & & & I & T \end{bmatrix} \quad (13.5)$$

where  $I$  is the  $n-1 \times n-1$  identity matrix, and  $T$  is the  $n-1 \times n-1$  tridiagonal matrix

$$\begin{bmatrix} -4 & 1 & & & \\ 1 & -4 & 1 & & \\ & \ddots & \ddots & \ddots & \\ & & 1 & -4 & 1 \\ & & & 1 & -4 \end{bmatrix}.$$

The vectors  $p$  and  $q$  come from the boundary conditions of (13.3), and are given by

$$p = \begin{bmatrix} p^1 \\ \vdots \\ p^{n-1} \end{bmatrix}, \quad q = \begin{bmatrix} q^1 \\ \vdots \\ q^{n-1} \end{bmatrix},$$

where

$$p^j = \frac{1}{h^2} \begin{bmatrix} g(x_0, y_j) \\ 0 \\ \vdots \\ 0 \\ g(x_n, y_j) \end{bmatrix}, \quad 1 \leq j \leq n-1,$$

and

$$q^1 = \frac{1}{h^2} \begin{bmatrix} g(x_1, y_0) \\ g(x_2, y_0) \\ \vdots \\ g(x_{n-2}, y_0) \\ g(x_{n-1}, y_0) \end{bmatrix}, \quad q^{n-1} = \frac{1}{h^2} \begin{bmatrix} g(x_1, y_n) \\ g(x_2, y_n) \\ \vdots \\ g(x_{n-2}, y_n) \\ g(x_{n-1}, y_n) \end{bmatrix}, \quad q^j = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \quad 2 \leq j \leq n-2.$$

The vector  $f$  comes from the source term of (13.3), and is given by

$$f = \begin{bmatrix} f^1 \\ \vdots \\ f^{n-1} \end{bmatrix}, \quad \text{where } f^j = \begin{bmatrix} f(x_1, y_j) \\ f(x_2, y_j) \\ \vdots \\ f(x_{n-1}, y_j) \end{bmatrix}$$

Note that this linear system is very large ( $A$  has  $(n-1)^4$  entries) and very sparse (most of the entries in  $A$  are zero). Thus we will should make use of sparse matrix routines (such as those in `scipy.sparse` and `scipy.sparse.linalg`) in order to reduce the time and memory used in setting up and solving the linear system.

**Problem 1.** Complete the function `poisson_square` by implementing the finite difference method 13.4. Use `scipy.sparse.linalg.spsolve` to solve the linear system. Use your function to solve the boundary value problem:

$$\begin{aligned} \Delta u &= 0, \quad x \in [0, 1] \times [0, 1], \\ u(x, y) &= x^3, \quad (x, y) \in \partial([0, 1] \times [0, 1]). \end{aligned} \tag{13.6}$$

Use  $n = 100$  subintervals for both  $x$  and  $y$ . Plot the solution as a 3D surface.

## Poisson's equation and conservative forces

In physics Poisson's equation is used to describe the scalar potential of a conservative force. In general

$$\Delta V = -f$$

where  $V$  is the scalar potential of the force, or the potential energy a particle would have at that point, and  $f$  is a source term. Examples of conservative forces include Newton's Law of Gravity (where matter become the source term) and Coulomb's Law, which gives the force between two charge particles (where charge is the source term).

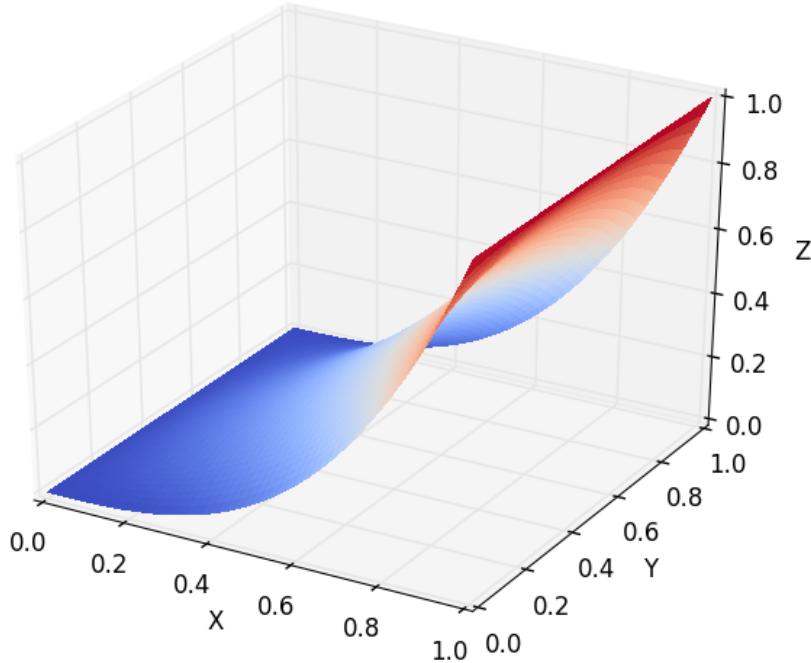


Figure 13.1: The solution of (13.6).

In electrostatics the electric potential is also known as the voltage, and is denoted by  $V$ . From Maxwell's equations it can be shown that that the voltage obeys Poisson's equation with the electric charge density (like a continuous cloud of electrons) being the source term:

$$\Delta V = -\frac{\rho}{\varepsilon_0},$$

where  $\rho$  is the charge density and  $\varepsilon_0$  is the permittivity of free space, which is a constant that we'll leave as 1.

Usually a non zero  $V$  at a point will cause a charged particle to move to a lower potential, changing  $\rho$  and the solution to  $V$ . However, in this analysis we'll assume that the charges are fixed in place.

Suppose we have 3 nested pipes. The outer pipe is attached to "ground," which usually we define to be  $V = 0$ , and the inner two have opposite relative charges. Physically the two inner pipes would function like a capacitor.

The following code will plot the charge distribution of this setup.

```
import matplotlib.colors as mcolors

def source(X,Y):
    """
    This function generates a 2D grid of points (X, Y) and returns a corresponding
    array of charge densities rho. The charge density is zero outside the unit square
    [0,1] x [0,1]. Inside the square, the charge density is zero except for a
    triangular region in the upper-right quadrant where it is positive. The
    triangular region has vertices at (0.6, 0.0), (1.0, 0.0), and (1.0, 0.6).
    The charge density is highest at the vertex (1.0, 0.6) and decreases towards
    the base of the triangle.
    """
    rho = np.zeros_like(X)
    mask = ((X > 0.6) & (Y < 0.6)) | ((X > 1.0) & (Y < 0.0))
    rho[mask] = 1.0
    return rho
```

```

Takes arbitrary arrays of coordinates X and Y and returns an array of the ←
    same shape
representing the charge density of nested charged squares
"""
src = np.zeros(X.shape)
src[ np.logical_or(
    np.logical_and( np.logical_or(abs(X-1.5) < .1,abs(X+1.5) < .1) ,abs(Y) ←
        < 1.6),
    np.logical_and( np.logical_or(abs(Y-1.5) < .1,abs(Y+1.5) < .1) ,abs(X) ←
        < 1.6))] = 1
src[ np.logical_or(
    np.logical_and( np.logical_or(abs(X-0.9) < .1,abs(X+0.9) < .1) ,abs(Y) ←
        < 1.0),
    np.logical_and( np.logical_or(abs(Y-0.9) < .1,abs(Y+0.9) < .1) ,abs(X) ←
        < 1.0))] = -1
return src

#Generate a color dictionary for use with LinearSegmentedColormap
#that places red and blue at the min and max values of data
#and white when data is zero

def genDict(data):
    zero = 1/(1 - np.max(data)/np.min(data))
    cdict = {'red': [(0.0, 1.0, 1.0),
                    (zero, 1.0, 1.0),
                    (1.0, 0.0, 0.0)],
             'green': [(0.0, 0.0, 0.0),
                        (zero, 1.0, 1.0),
                        (1.0, 0.0, 0.0)],
             'blue': [(0.0, 0.0, 0.0),
                       (zero, 1.0, 1.0),
                       (1.0, 1.0, 1.0)]}
    return cdict

a1 = -2.
b1 = 2.
c1 = -2.
d1 = 2.
n = 100
X = np.linspace(a1,b1,n)
Y = np.linspace(c1,d1,n)
X,Y = np.meshgrid(X,Y)

plt.imshow(source(X,Y),cmap = mcolors.LinearSegmentedColormap('cmap', genDict(←
    source(X,Y))))
plt.colorbar(label="Relative Charge")
plt.show()

```

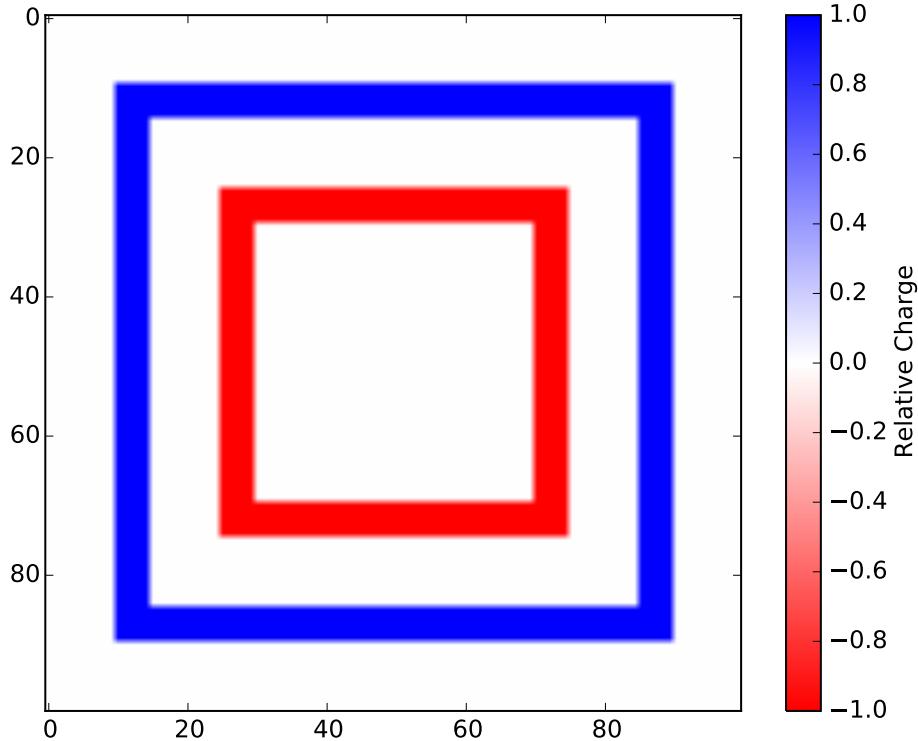


Figure 13.2: The charge density of the 3 nested pipes.

The function `genDict` scales the color values to be white when the charge density is zero. This is mostly to help visualize where there are neutrally charged zones by forcing them to be white. You may find it useful to also apply it when you solve for the electric potential.

With this definition of the charge density, we can solve Poisson's equation for the potential field.

**Problem 2.** Using the `poisson_square` function, solve

$$\begin{aligned} \Delta V &= -\rho(x, y), \quad x \in [-2, 2] \times [-2, 2], \\ u(x, y) &= 0, \quad (x, y) \in \partial([-2, 2] \times [-2, 2]). \end{aligned} \tag{13.7}$$

for the electric potential  $V$ . Use the source function defined above, such that  $\rho(x, y) = \text{source}(x, y)$ . Use  $n = 100$  subintervals for  $x$  and  $y$ . Use the provided code to plot your solution.

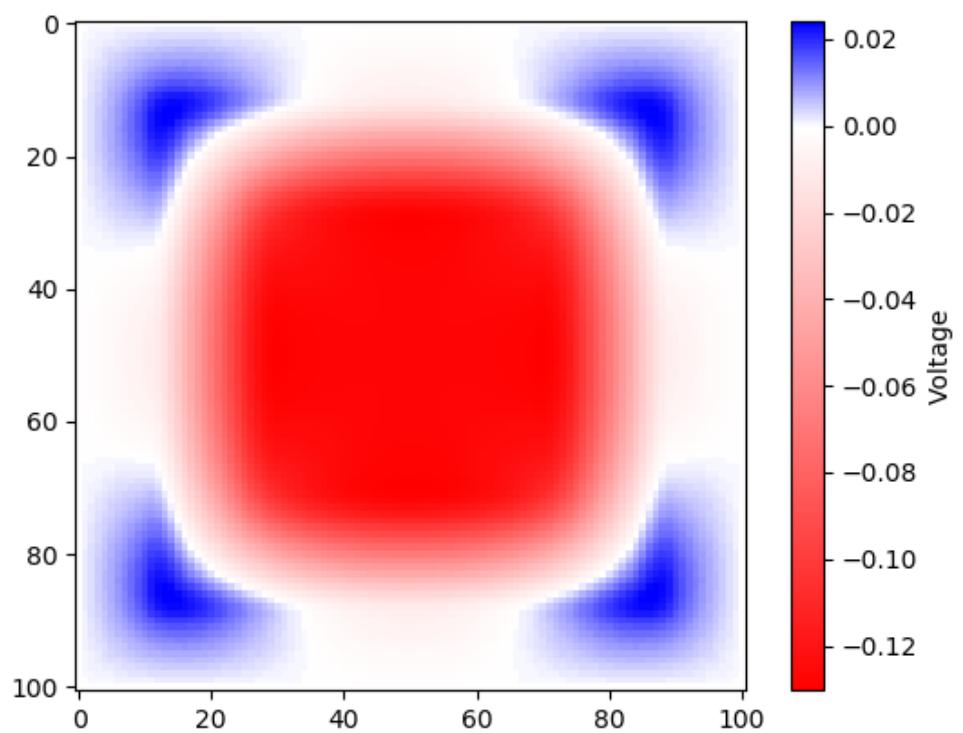


Figure 13.3: The electric potential of the 3 nested pipes.



# 14

## Method of Mean Weighted Residuals

**Lab Objective:** We introduce the method of mean weighted residuals (MWR) and use it to derive a pseudospectral method. This method will then be used to solve several boundary value problems.

Consider a linear differential equation

$$Lu = f,$$

defined on the interval  $[-1, 1]$ , together with associated boundary conditions. We will approximate the solution  $u(x)$  by a linear combination of  $N + 1$  basis functions  $\phi_i$ , so that

$$u(x) \approx u_N(x) = \sum_{i=0}^N a_i \phi_i(x).$$

To determine appropriate constants  $a_i$ , we then minimize the residual function

$$R(x, u_N) = Lu_N - f.$$

Note that  $R(x, u) = Lu - f = 0$  for the true solution  $u(x)$ .

This general strategy is often called the method of mean weighted residuals (MWR method). The MWR method is a general framework that describes many other, more specific methods. These more specific methods come from differing approaches to minimizing the residual  $R(x, u_N)$ , and the choice of basis functions  $\phi_i$ .

### The Pseudospectral Method

The pseudospectral or collocation method is obtained from the MWR method by forcing the residual function  $R(x, u_N)$  to equal zero at  $N + 1$  points in  $[-1, 1]$ , called collocation points. When done correctly, the pseudospectral method gives high accuracy and converges rapidly.

We will let the basis functions  $\phi_i$  be the Chebyshev polynomials,

$$T_0(x) = 1$$

$$T_1(x) = x$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x)$$

and the collocation points will be the Gauss-Lobatto points,  $x_i = \cos(\pi i/N)$ ,  $i = 0, \dots, N$ . The appropriate solution  $u_N$  may be represented with two equivalent forms. First,  $u_N$  can be described with the first  $N + 1$  coefficients  $\{a_i\}_{i=0}^N$  of its expansion in the Chebyshev polynomials. Since  $u_N$  is a polynomial of order  $N$ , it may be uniquely described by its values at the collocation points, that is, the unknown values  $\{u_N(x_i)\}_{i=0}^N$ .

These equivalent forms satisfy

$$MA = F \quad (14.1)$$

and

$$LU = F \quad (14.2)$$

where

$$\begin{aligned} U_i &= u(x_i), \\ A_i &= a_i, \\ F_i &= f(x_i), \\ L_{ij} &= (LC_j(x))|_{x=x_i}, \\ M_{ij} &= (L\phi_j(x))|_{x=x_i}. \end{aligned}$$

The functions  $C_j$  above are the cardinal functions, defined to be the polynomials of least degree satisfying

$$C_j(x_i) = \begin{cases} 1 & i = j \\ 0 & i \neq j. \end{cases}$$

Thus,  $u_N$  can also be expanded in the basis of the cardinal functions:

$$u_N(x) = \sum_{j=0}^N u_N(x_j)C_j(x).$$

When  $L = d/dx$ , the matrix corresponding to equation (14.2) is given by

$$L_{ij} = \frac{dC_j}{dx}(x_i) = \begin{cases} (1+2N^2)/6 & i = j = 0, \\ -(1+2N^2)/6 & i = j = N, \\ -x_j/[2(1-x_j^2)] & i = j, 0 < j < N, \\ (-1)^{i+j}\alpha_i/[\alpha_j(x_i - x_j)] & i \neq j. \end{cases}$$

where  $\alpha_0 = \alpha_N = 2$ , and  $\alpha_j = 1$  otherwise.

This matrix is often called the differentiation matrix ( $D$ ), and can be used to piece together the matrix  $L$  for more complicated differential operators. A stable, vectorized function to build the differentiation matrix is given below.

```
import numpy as np

def cheb(N):
    x = np.cos((np.pi/N)*np.linspace(0,N,N+1))
    x.shape = (N+1,1)
    lin = np.linspace(0,N,N+1)
```

```

lin.shape = (N+1,1)

c = np.ones((N+1,1))
c[0], c[-1] = 2., 2.
c = c*(-1.)**lin
X = x*np.ones(N+1) # broadcast along 2nd dimension (columns)

dX = X - X.T

D = (c*(1./c).T)/(dX + np.eye(N+1))
D = D - np.diag(np.sum(D.T, axis=0))
x.shape = (N+1,)
# Here we return the differentiation matrix and the Chebyshev points,
# numbered from x_0 = 1 to x_N = -1
return D, x

```

## Using the Differentiation Matrix

**Problem 1.** Use the differentiation matrix to numerically approximate the derivative of  $u(x) = e^x \cos(6x)$  on a grid of  $N$  Chebychev points where  $N = 6, 8$ , and  $10$ . (Use the linear system  $DU \approx U'$ .) Then use barycentric interpolation (`scipy.interpolate.barycentric_interpolate`) to approximate  $u'$  on a grid of  $100$  evenly spaced points.

Graphically compare your approximation to the exact derivative. Note that this convergence would not be occurring if the collocation points were equally spaced.

To approximate  $u''(x)$  on the grid  $\{x_i\}$ , we use

$$U'' \approx D^2U.$$

The BVP

$$\begin{aligned} u'' &= f(x), \quad x \in [-1, 1], \\ u(-1) &= 0, \quad u(1) = 0, \end{aligned}$$

can be discretized by the linear system

$$D^2U = F, \tag{14.3}$$

where  $F = [f(x_0), \dots, f(x_N)]^T$ . Since we have Dirichlet boundary conditions of 0, we can satisfy the boundary condition by forcing  $U[0] = U[N] = 0$ . This is done by replacing the first and last equations in (14.3) by the boundary conditions.

```

#The following code will force U[0] = U[N] = 0
D, x = cheb(N)      #for some N
D2 = np.dot(D, D)
D2[0,:], D2[-1,:] = 0, 0
D2[0,0], D2[-1,-1] = 1, 1
F[0], F[-1] = 0, 0

```

**Problem 2.** Use the pseudospectral method to solve the boundary value problem

$$\begin{aligned} u'' &= e^{2x}, \quad x \in (-1, 1), \\ u(-1) &= 0, \quad u(1) = 0. \end{aligned}$$

Use  $N = 8$  in the `cheb(N)` method and use barycentric interpolation to approximate  $u$  on 100 evenly spaced points. Compare your numerical solution with the exact solution,

$$u(x) = \frac{-\cosh(2) - \sinh(2)x + e^{2x}}{4}.$$

**Problem 3.** Use the pseudospectral method to solve the boundary value problem

$$\begin{aligned} u'' + u' &= e^{3x}, \quad x \in (-1, 1), \\ u(-1) &= 2, \quad u(1) = -1. \end{aligned}$$

Use  $N = 8$  in the `cheb(N)` method and use barycentric interpolation to approximate  $u$  on 100 evenly spaced points.

The previous exercise involved setting up and solving a linear system

$$AU = F,$$

where  $F$  is a vector whose entries are  $e^{3x}$  evaluated at the collocation points  $x_j$ , and  $U$  represents the approximation to the solution  $u$  at those points. However, whenever the ODE is nonlinear, the discretization becomes a nonlinear system of equations that must be solved using Newton's method. The next exercise contains a BVP whose ODE is nonlinear, with the additional complexity that the domain of the problem is not  $[-1, 1]$ .

**Problem 4.** Use the pseudospectral method to solve the boundary value problem

$$\begin{aligned} u'' &= \lambda \sinh(\lambda u), \quad x \in (0, 1), \\ u(0) &= 0, \quad u(1) = 1 \end{aligned}$$

for several values of  $\lambda$ :  $\lambda = 4, 8, 12$ . Begin by transforming this BVP onto the domain  $-1 < x < 1$ . Use  $N = 20$  in the `cheb(N)` method and use barycentric interpolation to approximate  $u$  on 100 evenly spaced points.

Below is sample code for implementing Newton's Method

```
from scipy.optimize import root

N = 20
D, x = cheb(20)

def F(U):
    out = None #Set up the equation you want the root of.
```

```
#Make sure to set the boundaries correctly

return out #Newtons Method will update U until the output is all 0's.

guess = None #Make your guess, same size as the cheb(N) output
solution = root(F, guess).x
```

## Minimizing the Area of a Surface of Revolution

A surface of revolution that minimizes its area is an example of a larger class of surfaces called minimal surfaces. A famous example of a minimal surface is a soap bubble. Soap bubbles minimize their surface area while containing a fixed volume of air. This behavior extends to merged bubbles, and a soap film whose boundary is a wire frame. Minimal surfaces have applications in molecular engineering and material science, and general relativity, where they describe the apparent horizon of a black hole.

Consider a function  $y(x)$  defined on  $[-1, 1]$  satisfying  $y(-1) = a$ ,  $y(1) = b$ . The area of the surface obtained by revolving the graph of  $y(x)$  about the  $x$ -axis is given by

$$T[y(x)] = \int_{-1}^1 2\pi y(x) \sqrt{1 + (y'(x))^2} dx.$$

To find the function  $y(x)$  whose surface of revolution minimizes surface area, we must minimize the functional  $T[y]$ . This is a classical problem from a branch of mathematics called the calculus of variations. Standard derivatives allow us to find the minimum values of functions defined on  $\mathbb{R}^n$ , and where they occur. The calculus of variations allows us to find the minimum values of functions whose input are other functions.

From the calculus of variations we know that a necessary condition for  $y(x)$  to minimize  $T[y]$  is that the Euler-Lagrange equation must be satisfied:

$$L_y - \frac{d}{dx} L_{y'} = 0,$$

where  $L(x, y, y') = 2\pi y \sqrt{1 + (y')^2}$ . Simplifying the Euler-Lagrange equation for our problem results in the ODE

$$yy'' - (y')^2 - 1 = 0.$$

Discretizing this ODE using the pseudospectral method results in the (nonlinear) system of equations

$$Y \cdot (D^2 Y) - (DY) \cdot (DY) = I,$$

where  $I$  is a vector of ones.

**Problem 5.** Find the function  $y(x)$  that satisfies  $y(-1) = 1$ ,  $y(1) = 7$ , and whose surface of revolution (about the  $x$ -axis) minimizes surface area. Compute the surface area, and plot the surface. Use  $N = 50$  in the `cheb(N)` method and use barycentric interpolation to approximate  $u$  on 100 evenly spaced points.

Below is sample code for creating the 3D wireframe figure.

```
from mpl_toolkits.mplot3d import Axes3D
```

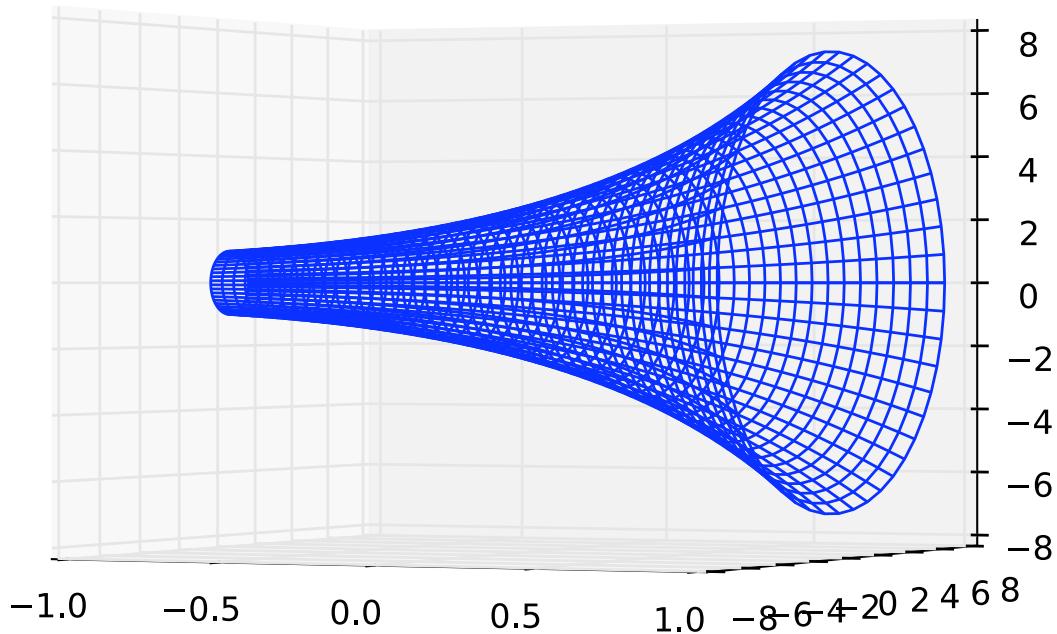


Figure 14.1: The minimal surface corresponding to Problem 5.

```
barycentric = None #This is the output of barycentric_interpolate() on ←
                   100 points

lin = np.linspace(-1, 1, 100)
theta = np.linspace(0, 2*np.pi, 401)
X, T = np.meshgrid(lin, theta)
Y, Z = barycentric*np.cos(T), barycentric*np.sin(T)

fig = plt.figure()
ax = fig.gca(projection="3d")
ax.plot_wireframe(X, Y, Z, rstride=10, cstride=10)
plt.show()
```

# 15

## A Pseudospectral method for periodic functions

**Lab Objective:** We look at a pseudospectral method with a Fourier basis, and numerically solve the advection equation using a pseudospectral discretization in space and a Runge-Kutta integration scheme in time.

Let  $f$  be a periodic function on  $[0, 2\pi]$ . Let  $x_1, \dots, x_N$  be  $N$  evenly spaced grid points on  $[0, 2\pi]$ . Since  $f$  is periodic on  $[0, 2\pi]$ , we can ignore the grid point  $x_N = 2\pi$ . We will further assume that  $N$  is even; similar formulas can be derived for  $N$  odd. Let  $h = 2\pi/N$ ; then  $\{x_0, \dots, x_{N-1}\} = \{0, h, 2h, \dots, 2\pi - h\}$ .

The discrete Fourier transform (DFT) of  $f$ , denoted by  $\hat{f}$  or  $\mathcal{F}(f)$ , is given by

$$\hat{f}(k) = h \sum_{j=0}^{N-1} e^{-ikx_j} f(x_j) \quad \text{where } k = -N/2 + 1, \dots, 0, 1, \dots, N/2.$$

The inverse DFT is then given by

$$f(x_j) = \frac{1}{2\pi} \sum_{k=-N/2}^{N/2} \frac{e^{ikx_j}}{c_k} \hat{f}(k), \quad j = 0, \dots, N-1, \quad (15.1)$$

where

$$c_k = \begin{cases} 2 & \text{if } k = -N/2 \text{ or } k = N/2, \\ 1 & \text{otherwise.} \end{cases} \quad (15.2)$$

The inverse DFT can then be used to define a natural interpolant (sometimes called a band-limited interpolant) by evaluating (15.1) at any  $x$  rather than  $x_j$ :

$$p(x) = \frac{1}{2\pi} \sum_{k=-N/2}^{N/2} e^{ikx} \hat{f}(k). \quad (15.3)$$

The interpolant for  $f'$  is then given by

$$p'(x) = \frac{1}{2\pi} \sum_{k=-N/2+1}^{N/2-1} ike^{ikx} \hat{f}(k). \quad (15.4)$$

Consider the function  $u(x) = \sin^2(x)\cos(x) + e^{2\sin(x+1)}$ . Using (15.4), the derivative  $u'$  may be approximated with the following code.<sup>1</sup> We note that although we only approximate  $u'$  at the Fourier grid points, (15.4) provides an analytic approximation of  $u'$  in the form of a trigonometric polynomial.

```

import numpy as np
from scipy.fftpack import fft, ifft
import matplotlib.pyplot as plt

N=24
x1 = (2.*np.pi/N)*np.arange(1,N+1)
f = np.sin(x1)**2.*np.cos(x1) + np.exp(2.*np.sin(x1+1))

# This array is reordered in Python to
# accomodate the ordering inside the fft function in scipy.
k = np.concatenate(( np.arange(0,N/2) ,
                     np.array([0]) , # Because hat{f}'(k) at k = N/2 is zero.
                     np.arange(-N/2+1,0,1) ))

# Approximates the derivative using the pseudospectral method
f_hat = fft(f)
fp_hat = ((1j*k)*f_hat)
fp = np.real(ifft(fp_hat))

# Calculates the derivative analytically
x2 = np.linspace(0,2*np.pi,200)
derivative = (2.*np.sin(x2)*np.cos(x2)**2. -
              np.sin(x2)**3. +
              2*np.cos(x2+1)*np.exp(2*np.sin(x2+1)))
)

plt.plot(x2,derivative,'-k',linewidth=2.)
plt.plot(x1,fp,'*b')
plt.savefig('spectral2_derivative.pdf')
plt.show()

```

**Problem 1.** Consider again the function  $u(x) = \sin^2(x)\cos(x) + e^{2\sin(x+1)}$ . Create a function that approximates  $\frac{1}{2}u'' - u'$  on the Fourier grid points for  $N = 24$ .

## The advection equation

Recall that the advection equation is given by

$$u_t + cu_x = 0 \quad (15.5)$$

---

<sup>1</sup>See Spectral Methods in MATLAB by Lloyd N. Trefethen. Another good reference is Chebyshev and Fourier Spectral Methods by John P. Boyd.

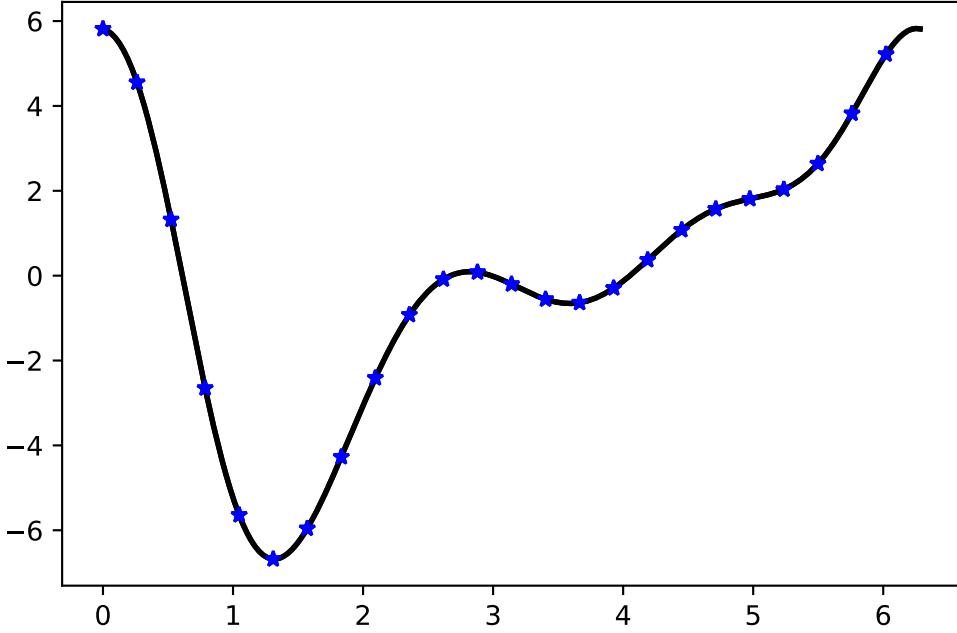


Figure 15.1: The derivative of  $u(x) = \sin^2(x) \cos(x) + e^{2 \sin(x+1)}$ .

where  $c$  is the speed of the wave (the wave travels to the right for  $c > 0$ ). We will consider the solution of the advection equation on the circle; this essentially amounts to solving the advection equation on  $[0, 2\pi]$  and assuming periodic boundary conditions.

A common method for solving time-dependent PDEs is called the method of lines. To apply the method of lines to our problem, we use our Fourier grid points in  $[0, \pi]$ : given an even  $N$ , let  $h = 2\pi/N$ , so that  $\{x_0, \dots, x_{N-1}\} = \{0, h, 2h, \dots, 2\pi - h\}$ . By using these grid points we obtain the collection of equations

$$u_t(x_j, t) + cu_x(x_j, t) = 0, \quad t > 0, \quad j = 0, \dots, N-1. \quad (15.6)$$

Let  $U(t)$  be the vector valued function given by  $U(t) = (u(x_j, t))_{j=0}^{N-1}$ . Let  $\mathcal{F}(U)(t)$  denote the discrete Fourier transform of  $u(x, t)$  (in space), so that

$$\mathcal{F}(U)(t) = (\hat{u}(k, t))_{k=-N/2+1}^{N/2}.$$

Define  $\mathcal{F}^{-1}$  similarly. Using the pseudospectral approximation in space leads to the system of ODEs

$$U_t + \vec{c} \mathcal{F}^{-1} \left( i \vec{k} \mathcal{F}(U) \right) = 0 \quad (15.7)$$

where  $\vec{k}$  is a vector, and  $\vec{k} \mathcal{F}(U)$  denotes element-wise multiplication. Similarly  $\vec{c}$  could also be a vector, if the wave speed  $c$  is allowed to vary.

**Problem 2.** Using a fourth order Runge-Kutta method (RK4), solve the initial value problem

$$u_t + c(x)u_x = 0, \quad (15.8)$$

where  $c(x) = .2 + \sin^2(x - 1)$ , and  $u(x, t = 0) = e^{-100(x-1)^2}$ . Plot your numerical solution from  $t = 0$  to  $t = 8$  over 150 time steps and 100  $x$  steps. Note that the initial data is nearly zero near  $x = 0$  and  $2\pi$ , and so we can use the pseudospectral method.<sup>a</sup> Use the following code to help graph.

```
t_steps = 150      # Time steps
x_steps = 100      # x steps

...
Your code here to set things up
...

sol = # RK4 method. Should return a t_steps by x_steps array

X,Y = np.meshgrid(x_domain, t_domain)
fig = plt.figure()
ax = fig.add_subplot(111, projection="3d")
ax.plot_wireframe(X,Y,sol)
ax.set_zlim(0,3)
plt.show()
```

---

<sup>a</sup>This problem is solved in Spectral Methods in MATLAB using a leapfrog discretization in time.

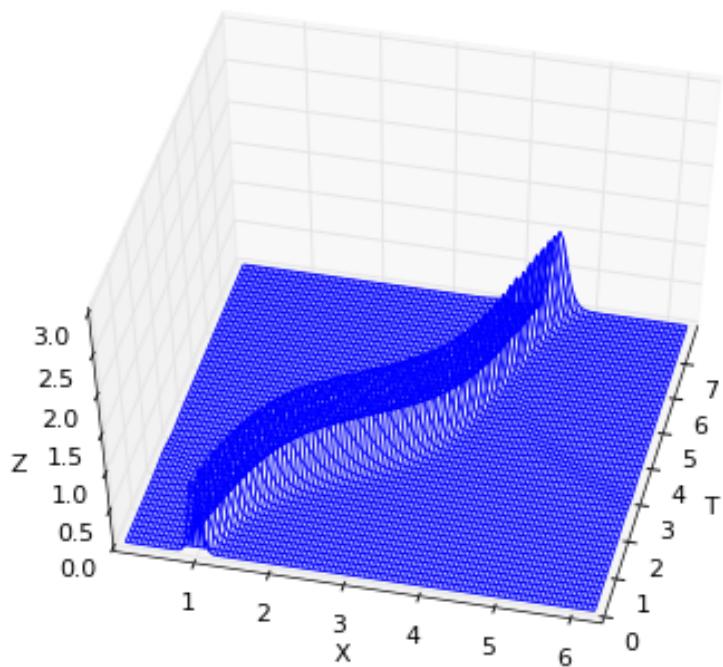


Figure 15.2: The solution of the variable speed advection equation; see Problem 2.



# 16 Inverse Problems

An important concept in mathematics is the idea of a well posed problem. The concept initially came from Jacques Hadamard. A mathematical problem is well posed if

1. a solution exists,
2. that solution is unique, and
3. the solution is continuously dependent on the data in the problem.

A problem that is not well posed is ill posed. Notice that a problem may be well posed, and yet still possess the property that small changes in the data result in larger changes in the solution; in this case the problem is said to be ill conditioned, and has a large condition number.

Note that for a physical phenomena, a well posed mathematical model would seem to be a necessary requirement! However, there are important examples of mathematical problems that are ill posed. For example, consider the process of differentiation. Given a function  $u$  together with its derivative  $u'$ , let  $\tilde{u}(t) = u(t) + \varepsilon \sin(\varepsilon^{-2}t)$  for some small  $\varepsilon > 0$ . Then note that

$$\|u - \tilde{u}\|_\infty = \varepsilon,$$

while

$$\|u' - \tilde{u}'\|_\infty = \varepsilon^{-1}.$$

Since a small change in the data leads to an arbitrarily large change in the output, differentiation is an ill posed problem. And we haven't even mentioned numerically approximating a derivative!

For an example of an ill posed problem from PDEs, consider the backwards heat equation with zero Dirichlet conditions:

$$\begin{aligned} u_t &= -u_{xx}, \quad (x, t) \in (0, L) \times (0, \infty), \\ u(0, t) &= u(L, t) = 0, \quad t \in (0, \infty), \\ u(x, 0) &= f(x), \quad x \in (0, L). \end{aligned} \tag{16.1}$$

For the initial data  $f(x)$ , the unique<sup>1</sup> solution is  $u(x, t) = 0$ . Given the initial data  $f(x) = \frac{1}{n} \sin(\frac{n\pi x}{L})$ , one can check that there is a unique solution  $u(x, t) = \frac{1}{n} \sin(\frac{n\pi x}{L}) \exp((\frac{n\pi}{L})^2 t)$ . Thus, on a finite interval  $[0, T]$ , as  $n \rightarrow \infty$  we see that a small difference in the initial data results in an arbitrarily large difference in the solution.

---

<sup>1</sup>See Partial Differential Equations by Lawrence C. Evans, chapter 2.3, for a proof of uniqueness.

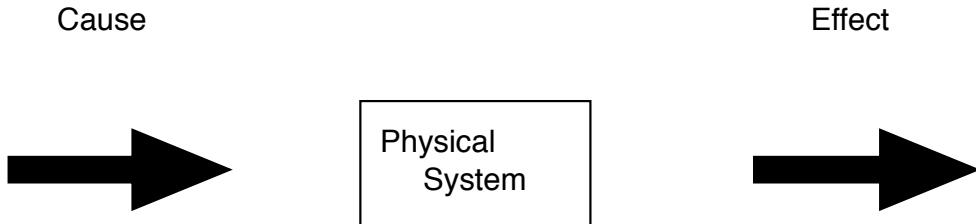


Figure 16.1: Cause and effect within a given physical system.

## Inverse Problems

As implied by the name, inverse problems come in pairs. For example, differentiation and integration are inverse problems. The easier problem (in this case integration) is often called the direct problem. Historically, the direct problem is usually studied first.

Given a physical system, together with initial data (the “cause”), the direct problem will usually predict the future state of the physical system (the “effect”); see Figure 16.1. Inverse problems often turn this on its head - given the current state of a physical system at time  $T$ , what was the physical state at time  $t = 0$ ?

Alternatively, suppose we measure the current state of the system, and we then measure the state at some future time. An important inverse problem is to determine an appropriate mathematical model that can describe the evolution of the system.

## Another look at heat flow through a rod

Consider the following ordinary differential equation, together with natural boundary conditions at the ends of the interval<sup>2</sup>:

$$\begin{cases} -(au')' = f, & x \in (0, 1), \\ a(0)u'(0) = c_0, & a(1)u'(1) = c_1. \end{cases} \quad (16.2)$$

This BVP can, for example, be used to describe the flow of heat through a rod. The boundary conditions would correspond to specifying the heat flux through the ends of the rod. The function  $f(x)$  would then represent external heat sources along the rod, and  $a(x)$  the density of the rod at each point.

---

<sup>2</sup>This example of an ill-posed problem is given in *Inverse Problems in the Mathematical Sciences* by Charles W Groetsch.

Typically, the density  $a(x)$  would be specified, along with any heat sources  $f(x)$ , and the (direct) problem is to solve for the steady-state heat distribution  $u(x)$ . Here we shake things up a bit: suppose the heat sources  $f$  are given, and we can measure the heat distribution  $u(x)$ . Can we find the density of the rod? This is an example of a parameter estimation problem.

Let us consider a numerical method for solving (16.2) for the density  $a(x)$ . Subdivide  $[0, 1]$  into  $N$  equal subintervals, and let  $x_j = jh$ ,  $j = 0, \dots, N$ , where  $h = 1/N$ . Let  $\phi_j(x)$  be the tent functions (used earlier in the finite element lab), given by

$$\phi_j(x) = \begin{cases} (x - x_{j-1})/h & x \in [x_{j-1}, x_j], \\ (x_{j+1} - x)/h & x \in [x_j, x_{j+1}], \\ 0 & \text{otherwise.} \end{cases}$$

We look for an approximation  $a^h(x)$  that is a linear combination of tent functions. This will be of the form

$$a^h = \sum_{j=0}^N \alpha_j \phi_j, \quad \alpha_j = a(x_j). \quad (16.3)$$

The  $h$  in this equation indicates that each of the tent functions in the linear combination rely on  $h = 1/N$ , and that a different  $h$  or  $N$  will result in different tent functions, so  $a^h$  will be different. The second half of (16.3) says that a good choice of  $a^h$  is found by taking  $\alpha_j = a(x_j)$ . Integrating (16.2) from 0 to  $x$ , we obtain

$$\begin{aligned} \int_0^x -(au')' ds &= \int_0^x f(s) ds, \\ -[a(x)u'(x) - c_0] &= \int_0^x f(s) ds, \\ u'(x) &= \frac{c_0 - \int_0^x f(s) ds}{a(x)}. \end{aligned} \quad (16.4)$$

Thus for each  $x_j$

$$\begin{aligned} u'(x_j) &= \frac{c_0 - \int_0^{x_j} f(s) ds}{a(x_j)}, \\ &= \frac{c_0 - \int_0^{x_j} f(s) ds}{\alpha_j}. \end{aligned}$$

The coefficients  $\alpha_j$  in (16.3) can now be approximated by minimizing

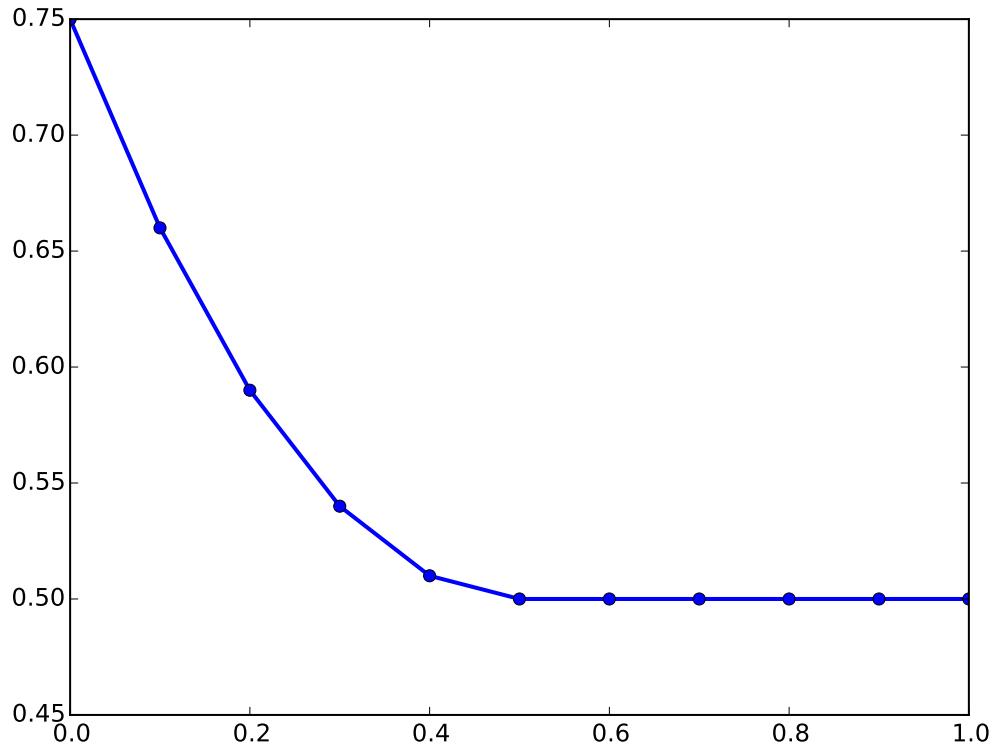
$$\left( \frac{c_0 - \int_0^{x_j} f(s) ds}{\alpha_j} - u'(x_j) \right)^2.$$

**Problem 1.** Solve (16.2) for  $a(x)$  using the following conditions:

$c_0 = 3/8$ ,  $c_1 = 5/4$ ,  $u(x) = x^2 + x/2 + 5/16$ ,  $x_j = .1j$  for  $j = 0, 1, \dots, 10$ , and

$$f = \begin{cases} -6x^2 + 3x - 1 & x \leq 1/2, \\ -1 & 1/2 < x \leq 1, \end{cases}$$

Produce the plot shown in Figure 16.2.

Figure 16.2: The solution  $a(x)$  to Problem 1

Hint: use the `minimize` function in `scipy.optimize` and some initial guess to find the  $a_j$ .

**Problem 2.** Find the density function  $a(x)$  satisfying

$$\begin{cases} -(au')' = -1, & x \in (0, 1), \\ a(0)u'(0) = 1, & a(1)u'(1) = 2. \end{cases} \quad (16.5)$$

where  $u(x) = x + 1 + \varepsilon \sin(\varepsilon^{-2}x)$ . Using several values of  $\varepsilon > 0.66049142$ , plot the corresponding density  $a(x)$  for  $x$  in `np.linspace(0, 1, 11)` to demonstrate that the problem is ill-posed.

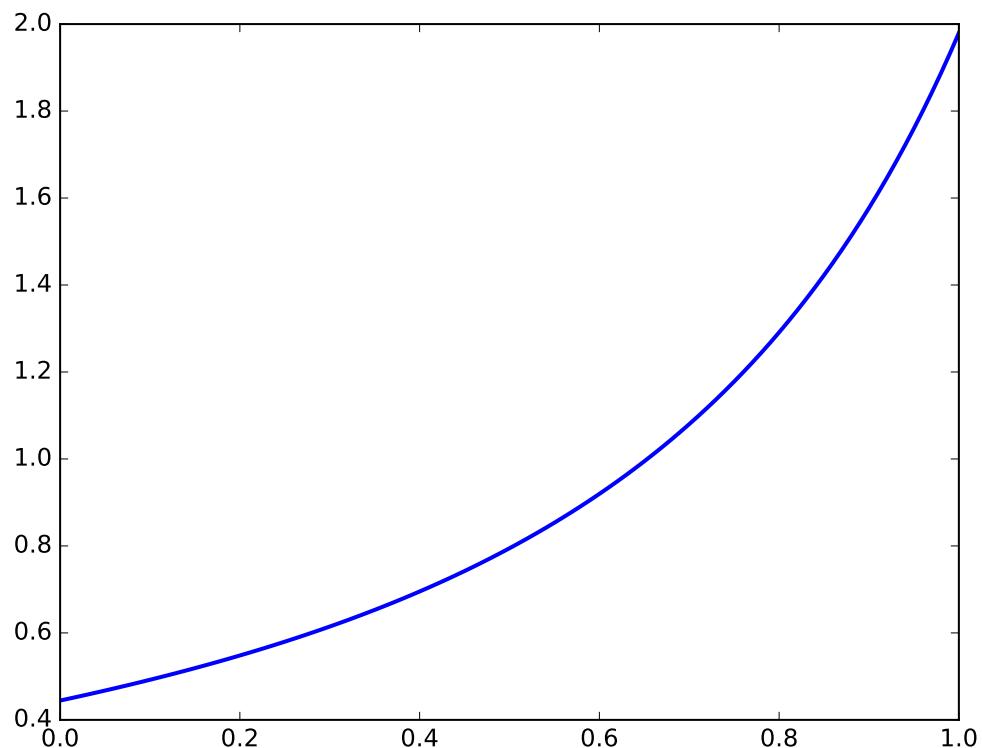


Figure 16.3: The density function  $a(x)$  satisfying (16.5) for  $\varepsilon = .8$ .



# 17

## The Shooting Method for Boundary Value Problems

Consider a boundary value problem of the form

$$\begin{aligned} y'' &= f(x, y, y'), \quad a \leq x \leq b, \\ y(a) &= \alpha, \quad y(b) = \beta. \end{aligned} \tag{17.1}$$

One natural way to approach this problem is to study the initial value problem (IVP) associated with this differential equation:

$$\begin{aligned} y'' &= f(x, y, y'), \quad a \leq x \leq b, \\ y(a) &= \alpha, \quad y'(a) = s. \end{aligned} \tag{17.2}$$

The goal is to determine an appropriate value  $s$  so that the solution of the IVP 17.2 is also a solution of the BVP 17.1.

Note that we can consider the value  $y(x)$  for a given initial condition  $y'(a) = s$  as a function of both  $x$  and  $s$ . Let  $y(x, s)$  be the solution of (17.2). The initial value conditions are then

$$y(a, s) = \alpha, \quad \frac{\partial y}{\partial x}(a, s) = s. \tag{17.3}$$

We wish to find a value of  $s$  so that  $y(b, s) = \beta$ . Consider the function  $h(s) = y(b, s) - \beta$ ; this function is called the residual function. If  $h(s) = 0$ , then  $y(b, s) = \beta$  and the boundary condition is satisfied, so zeros of the function  $h$  correspond to initial conditions that are solutions to the BVP 17.1. Applying Newton's method to the function  $h(s)$ , we obtain the iterative method

$$\begin{aligned} s_{n+1} &= s_n - \frac{h(s_n)}{h'(s_n)}, \\ &= s_n - \frac{y(b, s_n) - \beta}{\frac{\partial}{\partial s} y(b, s)|_{s_n}}, \quad n = 0, 1, \dots \end{aligned}$$

Provided our initial guess  $s_0$  is sufficiently good, this will converge to a value of  $s$  such that the initial value problem is also a solution to the boundary value problem. Notice that finding  $y(b, s_n)$  requires solving the initial value problem using RK4 or some other method.

We recall that Newton's method generally requires a good initial guess  $s_0$ . A plausible initial guess for this setup would be the average rate of change of the solution across the entire interval, which gives  $s_0 = (\beta - \alpha)/(b - a)$ . If this initial guess is insufficient, it may be refined by manually inspecting the solution  $y(x, s_0)$  of the initial value problem.

Using Newton's method requires us to evaluate or approximate the function  $h'(s_n)$ . This term may be approximated with a finite difference  $h'(s_n) \approx \frac{h(s_n) - h(s_{n-1})}{s_n - s_{n-1}}$ , giving us the iterative method

$$\begin{aligned}s_{n+1} &= s_n - h(s_n) \frac{(s_n - s_{n-1})}{h(s_n) - h(s_{n-1})} \\&= s_n - (y(b, s_n) - \beta) \frac{s_n - s_{n-1}}{y(b, s_n) - y(b, s_{n-1})}, \quad n = 1, 2, \dots\end{aligned}$$

This variation of Newton's method is called the secant method, and requires two initial values instead of one. The secant method generally does not converge as quickly as standard Newton's method, but it avoids needing to compute the actual derivative of  $h(s)$ .

As an example, consider the boundary value problem

$$\begin{aligned}y'' &= -4y - 9 \sin(x), \quad x \in [0, 3\pi/4], \\y(0) &= 1, \\y(3\pi/4) &= -\frac{1 + 3\sqrt{2}}{2}.\end{aligned}\tag{17.4}$$

This has the exact solution

$$y(x) = \cos(2x) + \frac{1}{2} \sin(2x) - 3 \sin(x).$$

The following code implements the secant method to solve (17.4) numerically. We use `scipy.integrate.solve_ivp` to solve the initial value problems.

```
import numpy as np
from scipy.integrate import solve_ivp
from matplotlib import pyplot as plt

# Secant method
def secant_method(h, s0, s1, max_iter=100, tol=1e-8):
    """
    Finds a root of h(s)=0 using the secant method with the
    initial guesses s0, s1.
    """
    for i in range(max_iter):
        # Get the residuals
        h0 = h(s0)
        h1 = h(s1)
        # Update
        s2 = s1 - h1 * (s1 - s0) / (h1 - h0)
        s0, s1 = s1, s2

        # Check convergence
        if abs(h1) < tol:
            return s2

    print("Secant method did not converge")
    return s2
```

```
# Define the ODE right-hand side
def ode(x, y):
    return np.array([
        y[1],
        -4*y[0]-9*np.sin(x)
    ])

# Endpoint values
a = 0
b = 3/4 * np.pi
alpha = 1
beta = - (1+3*np.sqrt(2))/2

# Define a residual function
def residual(s):
    # Find the right endpoint
    sol = solve_ivp(ode, (a, b), [alpha, s])
    yb = sol.y[0,-1]
    return yb - beta

# Find the right value of s using the secant method
s = secant_method(residual, (beta-alpha)/2, -1)

# Compute and plot the solution
x = np.linspace(0,3*np.pi/4,100)
y = solve_ivp(ode, (a,b), (alpha, s), t_eval=x).y[0]

plt.plot(x, y)
plt.show()
```

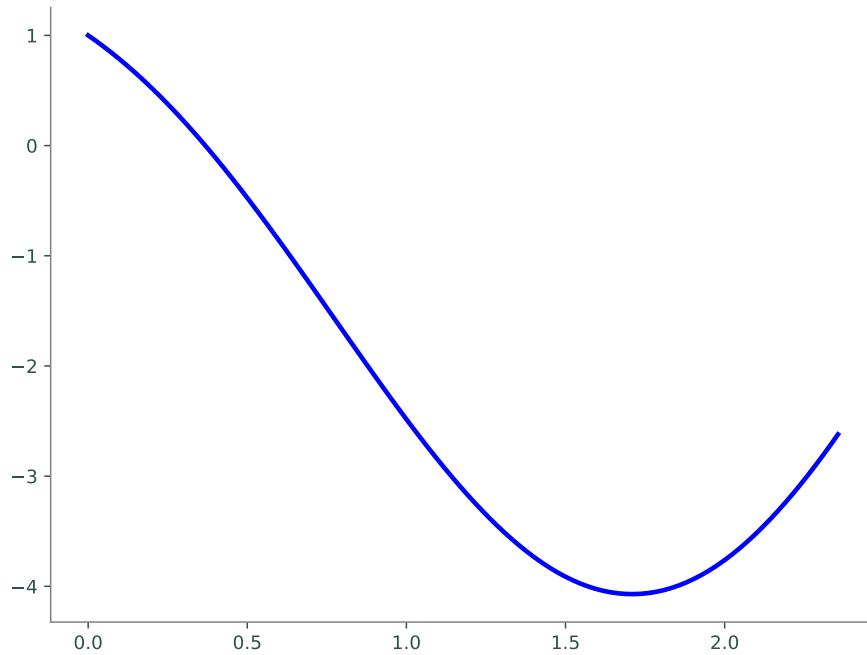


Figure 17.1: The solution to the BVP (17.4) from the above example.

**Problem 1.** Appropriately defined initial value problems will usually have a unique solution. Boundary value problems are not so straightforward; they may have no solution or they may have several, and you may have to determine which solutions are physically interesting.

Use the secant method to solve the following BVP:<sup>a</sup>

$$\begin{aligned} y'' &= -e^{y-1}, \quad x \in [0, 1], \\ y(0) &= y(1) = 1. \end{aligned}$$

This BVP has two solutions. Using the secant method, find both numerical solutions and their initial slopes. (Their plots are given in Figure 17.2.) What initial values  $s_0, s_1$  did you use to find them?

---

<sup>a</sup>This example is from Numerical Solution of Boundary Value Problems for Ordinary Differential Equations by Ascher, Mattheij, and Russell, page 89.

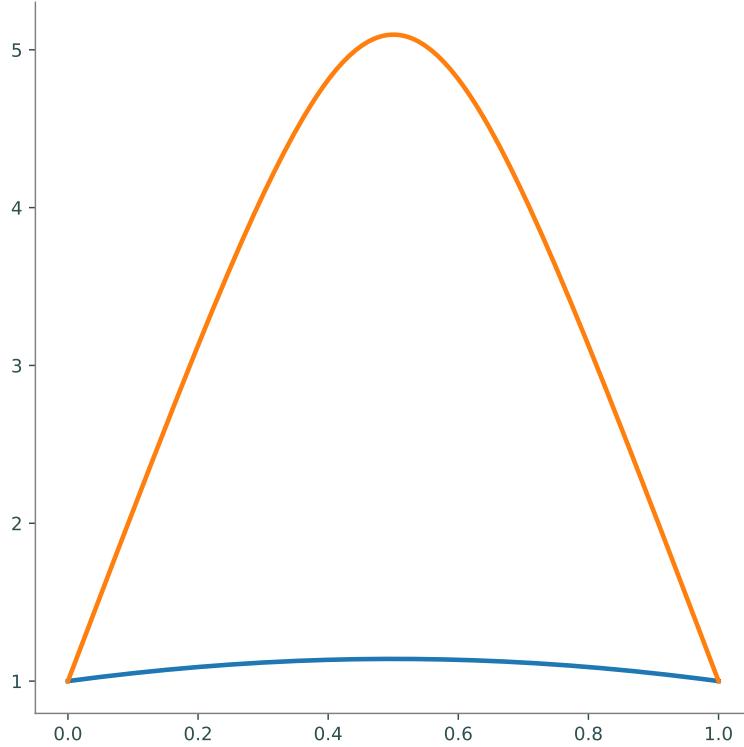


Figure 17.2: Both solutions to the boundary value problem given in Problem 1.

Instead of using the secant method, let us consider how to solve for  $h'(s) = \frac{\partial}{\partial s}y(b, s)$  for problems of the form given in (17.1). For typical systems of ODEs, the solution  $y(x, s)$  is smooth enough that it can be differentiated with respect to  $x$  and  $s$  in any order.<sup>1</sup> Let  $z(x, s) = \frac{\partial}{\partial s}y(x, s)$ , and note that  $h'(s) = z(b, s)$ . Using the chain rule, we obtain

$$\begin{aligned} z'' &= \frac{\partial}{\partial s}y''(x, s) = \frac{\partial f}{\partial y}(x, y(x, s), y'(x, s)) \cdot \frac{dy}{ds}(x, s), \\ &\quad + \frac{\partial f}{\partial y'}(x, y(x, s), y'(x, s)) \cdot \frac{\partial y'}{\partial s}(x, s), \end{aligned}$$

Using the initial conditions associated with  $y(x, s)$  and noting that  $z(x, s) = \frac{\partial}{\partial s}y(x, s)$  and  $z'(x, s) = \frac{\partial}{\partial s}y'(x, s)$ , we obtain the following initial value problem for  $z(x, s)$ :

$$\begin{aligned} z'' &= z \frac{\partial f}{\partial y}(x, y, y') + z' \frac{\partial f}{\partial y'}(x, y, y'), \quad a \leq x \leq b, \\ z(a, s) &= 0, \quad z'(a, s) = 1. \end{aligned}$$

---

<sup>1</sup>This is guaranteed to be the case if the right hand side of the ODE is  $C^1$ , as in all of the examples here, as this guarantees both partial derivatives of  $y$  are continuous.

To use Newton's method, the IVPs for  $y$  and  $z$  must be solved simultaneously. The iterative method then becomes

$$\begin{aligned}s_{n+1} &= s_n - \frac{h(s)}{h'(s)} \\&= s_n - \frac{y(b, s_n) - \beta}{z(b, s_n)}, \quad n = 0, 1, \dots\end{aligned}$$

**Problem 2.** Use Newton's method to solve the BVP

$$\begin{aligned}y'' &= 3 + \frac{2y}{x^2}, \quad x \in [1, e], \\y(1) &= 6, \\y(e) &= e^2 + 6/e.\end{aligned}$$

Plot your solution, and compare with Figure 17.3. What is an appropriate initial guess?

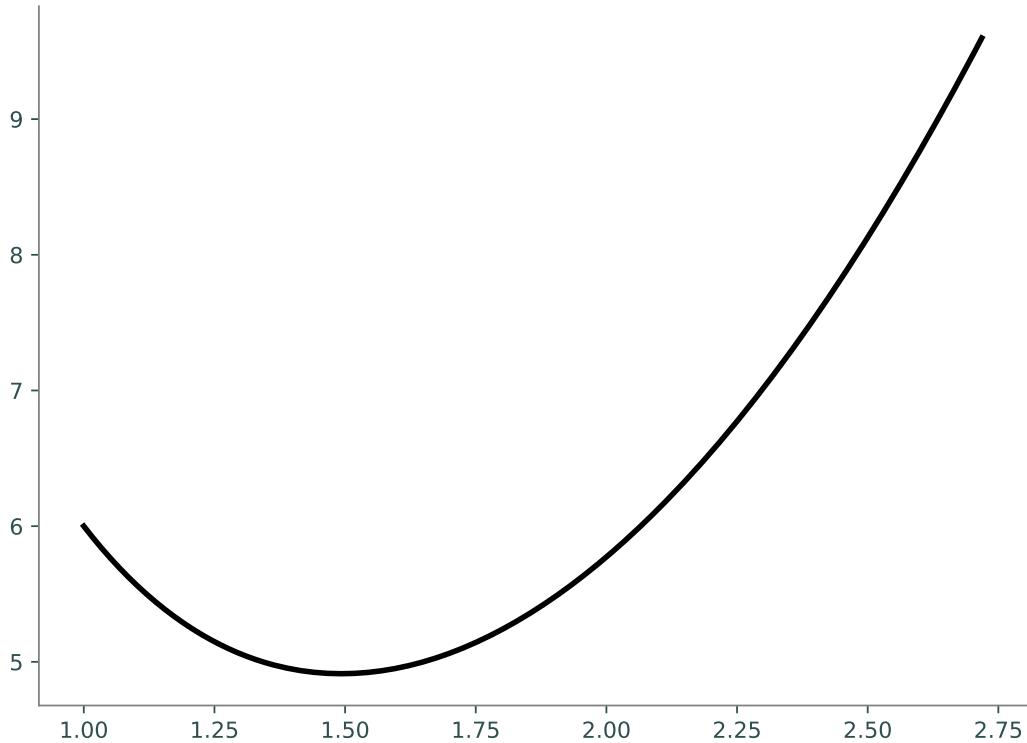


Figure 17.3: The solution of  $y'' = 3 + 2y/x^2$ , satisfying the boundary conditions  $y(1) = 6$ ,  $y(e) = e^2 + 6/e$ , given in Problem 2.

## The Cannon Problem

Consider the problem of aiming a projectile at a given target. Here we will construct a differential equation that describes the path of the projectile and takes into account air resistance. We will then use the shooting method to determine the angle at which the projectile should be launched.

Let  $t$  denote time, and the coordinates of the projectile be given by  $\mathbf{r}(t) = (x(t), y(t))$ . If  $\theta(t)$  represents the angle of the velocity vector from the positive  $x$ -axis and  $v(t) = \|\mathbf{v}(t)\|$  represents the speed of the projectile, then we have

$$\begin{aligned}\dot{x} &= v \cos \theta, \\ \dot{y} &= v \sin \theta.\end{aligned}$$

Note that each of  $x, y, \theta$ , and  $v$  are functions of  $t$ , so the dot denotes  $\frac{d}{dt}$ . The tangent vector to the path traced by the projectile is the unit vector in the direction of the projectile's velocity, so  $\mathbf{T}(t) = (\cos \theta, \sin \theta)$ . The unit normal vector  $\mathbf{N}(t)$  is given by  $\mathbf{N}(t) = (-\sin \theta, \cos \theta)$ . Thus the relationship between basis vectors  $\mathbf{i}, \mathbf{j}$ , and  $\mathbf{T}(t), \mathbf{N}(t)$  is given by

$$\begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} \mathbf{i} \\ \mathbf{j} \end{bmatrix} = \begin{bmatrix} \mathbf{T}(t) \\ \mathbf{N}(t) \end{bmatrix}$$

Let  $F_g$  represent the force on the projectile due to gravity, and  $F_d$  represent the force on the projectile due to air resistance. (We assume the air is still.) From Newton's law we have

$$m\dot{\mathbf{v}} = F_g + F_d.$$

The drag equation from fluid dynamics says that the force on the projectile due to air resistance is  $kv^2 = (1/2)\rho c_D A v^2$ , where  $\rho$  is the mass density of air (about  $1.225 \text{ kg/m}^3$ ),  $v$  is the speed of the projectile, and  $A$  is its cross-sectional area. The drag coefficient  $c_D$  is a dimensionless quantity that changes with respect to the shape of the object. (If we assume our projectile is spherical with a diameter of .2 m, then its drag coefficient  $c_D \approx 0.47$ , its cross-sectional area is  $\pi/100 \text{ m}^2$ , and we obtain  $k \approx 0.009$ .)

Thus the total force on the shell is

$$\begin{aligned}m\dot{\mathbf{v}} &= -mg\mathbf{j} - kv^2\mathbf{T}, \\ &= -mg(\mathbf{T} \sin \theta + \mathbf{N} \cos \theta) - kv^2\mathbf{T}, \\ &= (-mg \sin \theta - kv^2)\mathbf{T} - mg\mathbf{N} \cos \theta.\end{aligned}\tag{17.5}$$

From the identity  $\mathbf{v} = (\dot{x}, \dot{y}) = (v \cos \theta, v \sin \theta)$  we have

$$\begin{aligned}m\dot{\mathbf{v}} &= m(\dot{v} \cos \theta - v\dot{\theta} \sin \theta, \dot{v} \sin \theta + v\dot{\theta} \cos \theta) \\ &= m(\dot{v} \cos \theta - v\dot{\theta} \sin \theta)(\cos \theta \mathbf{T} - \sin \theta \mathbf{N}) \\ &\quad + m(\dot{v} \sin \theta + v\dot{\theta} \cos \theta)(\mathbf{T} \sin \theta + \mathbf{N} \cos \theta), \\ &= m(\mathbf{T}\dot{v} + \mathbf{N}v\dot{\theta}).\end{aligned}\tag{17.6}$$

From equations (17.5) and (17.6) we have

$$\begin{aligned}m\dot{v} &= -mg \sin \theta - kv^2, \\ m\dot{v}\dot{\theta} &= -mg \cos \theta.\end{aligned}$$

Thus we have the system of differential equations

$$\begin{aligned}\dot{x} &= v \cos \theta, \\ \dot{y} &= v \sin \theta, \\ \dot{v} &= -g \sin \theta - kv^2/m, \\ \dot{\theta} &= -g \cos \theta/v.\end{aligned}$$

We can actually write this problem to be independent of  $t$ , which will make solving it simpler, since we do not know the final time of impact. If we assume that  $t$  is an smooth invertible function of  $x$  (that is,  $t = t(x)$ ), then we obtain

$$\begin{aligned}\frac{dy}{dx} &= \frac{dy}{dt} \frac{dt}{dx}, \\ &= \frac{dy}{dt} \frac{1}{\frac{dx}{dt}}, \\ &= \frac{v \sin \theta}{v \cos \theta} = \tan \theta.\end{aligned}$$

We find  $\frac{dv}{dx}$  and  $\frac{d\theta}{dx}$  in a similar manner. Thus our system of differential equations becomes

$$\begin{aligned}\frac{dy}{dx} &= \tan \theta, \\ \frac{dv}{dx} &= -\frac{g \sin \theta + \mu v^2}{v \cos \theta}, \\ \frac{d\theta}{dx} &= -\frac{g}{v^2},\end{aligned}\tag{17.7}$$

where  $\mu = k/m$ . We can now consider  $y, v$ , and  $\theta$  to be functions of  $x$ , and  $x$  to be the independent variable.

**Problem 3.** Suppose we have a cannon that fires a projectile at a velocity of 45 m/s, and the projectile has a mass of about 60 kg, so that  $\mu = .0003$ . At what angle  $\theta(0)$  should it be fired to land at a distance of 195 m? Use the secant method to find initial values for  $\theta$  that give solutions to the following BVP:

$$\begin{aligned}\frac{dy}{dx} &= \tan \theta, \\ \frac{dv}{dx} &= -\frac{g \sin \theta + \mu v^2}{v \cos \theta}, \\ \frac{d\theta}{dx} &= -\frac{g}{v^2}, \\ y(0) &= y(195) = 0, \\ v(0) &= 45 \text{ m/s}\end{aligned}\tag{17.8}$$

$$(g = 9.8067 \text{ m/s}^2.)$$

There are four initial angles  $\theta(0)$  that produce solutions for this BVP when  $\mu = 0.0003$ . Find and plot at least two of them.<sup>a</sup> Also find the two solutions when  $\mu = 0$  (no air resistance), and compare. Graphs of the solutions are given in Figure 17.5.

Keep in mind that the unknown initial condition is  $\theta(0)$ , not  $y'(0)$ . What is the appropriate residual function  $h(t)$  to apply the secant method to?

<sup>a</sup>For two of them, it is difficult to get the secant method to converge to their initial values.

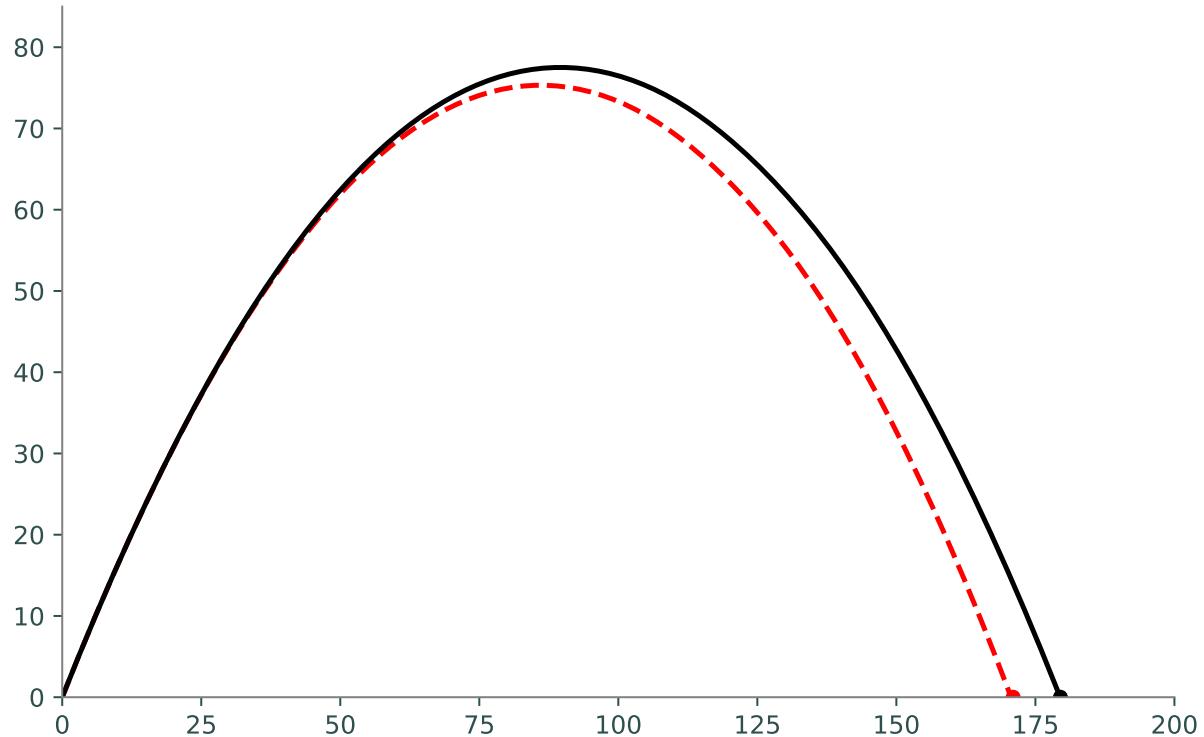


Figure 17.4: Two solutions of the system of equations (17.7), both with initial conditions  $y(0) = 0$  m,  $v(0) = 45$  m/s, and  $\theta(0) = \pi/3$ . The black curve is the trajectory of a projectile with no air resistance ( $\mu = 0$ ). The red curve describes the trajectory of a more realistic projectile ( $\mu = .0003$ ).

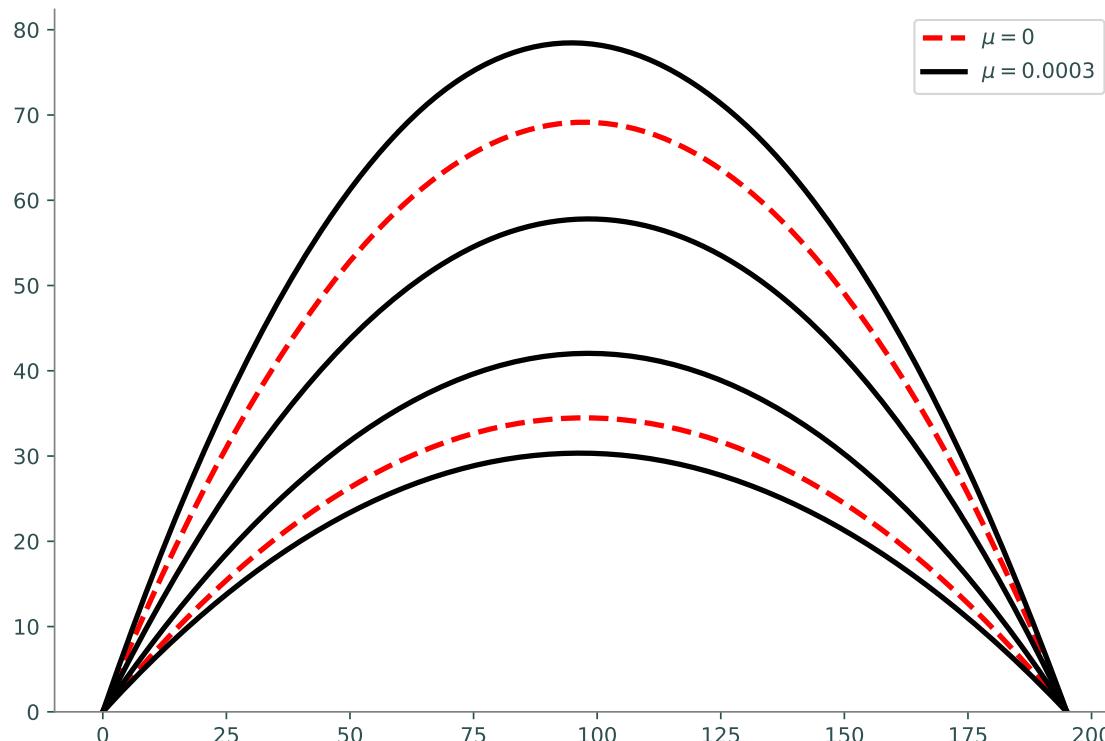


Figure 17.5: All four solutions of the boundary value problem (17.8) when the air resistance is described by the parameter  $\mu = .0003$ . Also shows both solutions with no air resistance ( $\mu = 0$ ).

# 18

## Total Variation and Image Processing

**Lab Objective:** Minimizing an energy functional is equivalent to solving the resulting Euler-Lagrange equations. We introduce the method of steepest descent to solve these equations, and apply this technique to a denoising problem in image processing.

### The Gradient Descent method

Consider an energy functional  $J[u]$ , defined over a collection of admissible functions  $u : \Omega \subset \mathbb{R}^n \rightarrow \mathbb{R}$ , with the form

$$J[u] = \int_{\Omega} L(x, u, \nabla u) dx$$

where  $L = L(x, u, \nabla u)$  is a function  $\mathbb{R}^n \times \mathbb{R} \times \mathbb{R}^n \rightarrow \mathbb{R}$ . A standard result from the calculus of variations states that a minimizing function  $u^*$  satisfies the Euler-Lagrange equation

$$L_u - \sum_{i=1}^n \frac{\partial L_{u_{x_i}}}{\partial x_i} = L_u - \nabla \cdot L_{\nabla u} = L_u - \operatorname{div}(L_{\nabla u}) = 0. \quad (18.1)$$

where  $L_{\nabla u} = \nabla' L = [L_{x_1}, \dots, L_{x_n}]^\top$ .

This equation is typically an elliptic PDE, possessing boundary conditions associated with restrictions on the class of admissible functions  $u$ . To more easily compute (18.1), we consider a related parabolic PDE,

$$\begin{aligned} u_t &= -(L_y - \operatorname{div} L_{\nabla u}), & t > 0, \\ u(x, 0) &= u_0(x), & t = 0. \end{aligned} \quad (18.2)$$

A steady state solution of (18.2) does not depend on time, and thus solves the Euler-Lagrange equation. It is often easier to evolve an initial guess using (18.2), and stop whenever its steady state is well-approximated, than to solve (18.1) directly.

Consider the energy functional

$$J[u] = \int_{\Omega} \|\nabla u\|^2 dx.$$

The minimizing function  $u^*$  satisfies the Euler-Lagrange equation

$$-\operatorname{div} \nabla u = -\Delta u = 0.$$

The gradient descent flow is the well-known heat equation

$$u_t = \Delta u.$$

The Euler-Lagrange equation could equivalently be described as  $\Delta u = 0$ , leading to the PDE  $u_t = -\Delta u$ . Since the backward heat equation is ill-posed, it would not be helpful in a search for the steady-state.

Let us take the time to make (18.2) more rigorous. We recall that

$$\begin{aligned} \delta J(u; h) &= \frac{d}{dt} J(u + \varepsilon h) \Big|_{\varepsilon=0}, \\ &= \int_{\Omega} (L_y(u) - \operatorname{div} L_{\nabla u}(u)) h \, dx, \\ &= \langle L_y(u) - \operatorname{div} L_{\nabla u}(u), h \rangle_{L^2(\Omega)}, \end{aligned}$$

for each  $u$  and each admissible perturbation  $h$ . Then using the Cauchy-Schwarz inequality,

$$|\delta J(u; h)| \leq \|L_y(u) - \operatorname{div} L_{\nabla u}(u)\| \cdot \|h\|$$

with equality iff  $h = \alpha(L_y(u) - \operatorname{div} L_{\nabla u}(u))$  for some  $\alpha \in \mathbb{R}$ . This implies that the “direction”  $h = L_y(u) - \operatorname{div} L_{\nabla u}(u)$  is the direction of steepest ascent and maximizes  $\delta J(u; h)$ . Similarly,

$$h = -(L_y(u) - \operatorname{div} L_{\nabla u}(u))$$

points in the direction of steepest descent, and the flow described by (18.2) tends to move toward a state of lesser energy.

## Minimizing the area of a surface of revolution

The area of the surface obtained by revolving a curve  $y(x)$  about the  $x$ -axis is

$$A[y] = \int_a^b 2\pi y \sqrt{1 + (y')^2} \, dx.$$

To minimize the functional  $A$  over the collection of smooth curves with fixed end points  $y(a) = y_a$ ,  $y(b) = y_b$ , we use the Euler-Lagrange equation

$$\begin{aligned} 0 &= 1 - y \frac{y''}{1 + (y')^2}, \\ &= 1 + (y')^2 - yy'', \end{aligned} \tag{18.3}$$

with the gradient descent flow given by

$$\begin{aligned} u_t &= -1 - (y')^2 + yy'', \quad t > 0, x \in (a, b), \\ u(x, 0) &= g(x), \quad t = 0, \\ u(a, t) &= y_a, \quad u(b, t) = y_b. \end{aligned} \tag{18.4}$$

## Numerical Implementation

We will construct a numerical solution of (18.4) using the conditions  $y(-1) = 1$ ,  $y(1) = 7$ . A simple solution can be found by using a second-order order discretization in space with a simple forward Euler step in time. We create the grid and set our end states below.

```

import numpy as np

a, b = -1, 1.
alpha, beta = 1., 7.
#### Define variables x_steps, final_T, time_steps ####
delta_t, delta_x = final_T/time_steps, (b-a)/x_steps
x0 = np.linspace(a,b,x_steps+1)

```

Most numerical schemes have a stability condition that must be satisfied. Our discretization requires that  $\frac{\Delta t}{(\Delta x)^2} \leq \frac{1}{2}$ . We continue by checking that this condition is satisfied, and use the straight line connecting the end points as initial data.

```

# Check a stability condition for this numerical method
if delta_t/delta_x**2. > .5:
    print("stability condition fails")

u = np.empty((2,x_steps+1))
u[0] = (beta - alpha)/(b-a)*(x0-a) + alpha
u[1] = (beta - alpha)/(b-a)*(x0-a) + alpha

```

Finally, we define the right hand side of our difference scheme, and time step until the scheme converges.

```

def rhs(y):
    # Approximate first and second derivatives to second order accuracy.
    yp = (np.roll(y,-1) - np.roll(y,1))/(2.*delta_x)
    ypp = (np.roll(y,-1) - 2.*y + np.roll(y,1))/delta_x**2.
    # Find approximation for the next time step, using a first order Euler step
    y[1:-1] -= delta_t*(1. + yp[1:-1]**2. - 1.*y[1:-1]*ypp[1:-1])

    # Time step until successive iterations are close
iteration = 0
while iteration < time_steps:
    rhs(u[1])
    if norm(np.abs((u[0] - u[1]))) < 1e-5: break
    u[0] = u[1]
    iteration+=1

print("Difference in iterations is ", norm(np.abs((u[0] - u[1]))))
print("Final time = ", iteration*delta_t)

```

**Problem 1.** Using 20  $x$  steps, 250 time steps,  $a = -1$ ,  $b = 1$ ,  $\alpha = 1$ ,  $\beta = 7$ , and a final time of 0.2, plot the solution that minimizes (18.4). It should match figure 18.1.

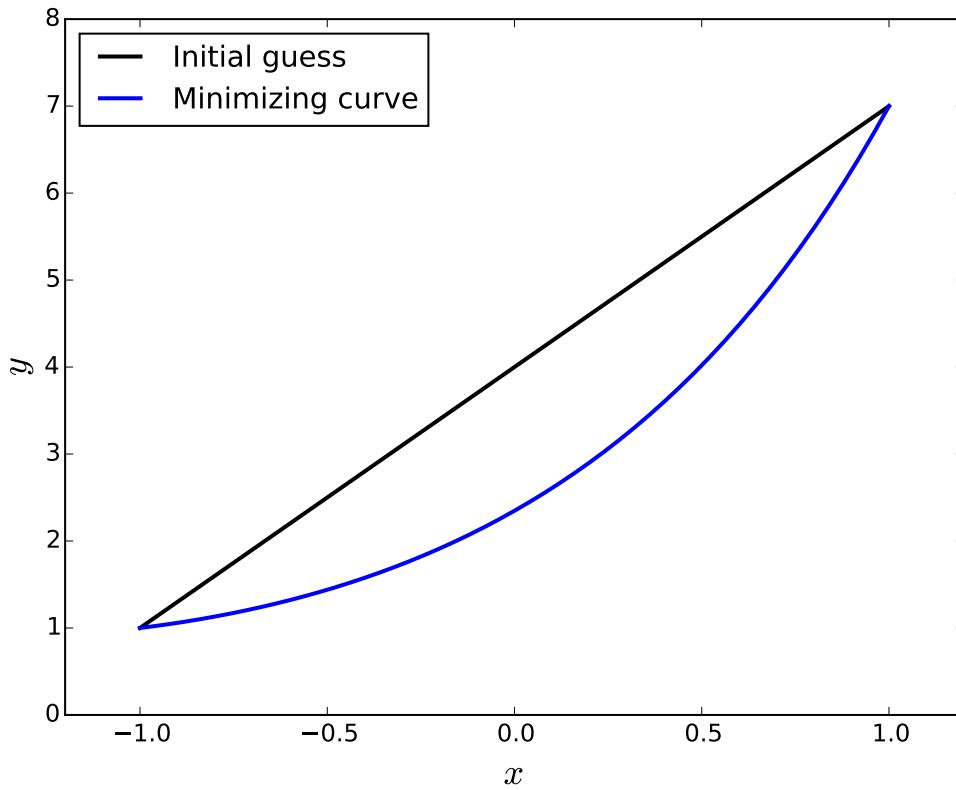


Figure 18.1: The solution of (18.3), found using the gradient descent flow (18.4).

## Image Processing: Denoising

A greyscale image can be represented by a scalar-valued function  $u : \Omega \rightarrow \mathbb{R}$ ,  $\Omega \subset \mathbb{R}^2$ . The following code reads an image into an array of floating point numbers, adds some noise, and saves the noisy image.

```
from numpy.random import randint, uniform, randn
import matplotlib.pyplot as plt
from matplotlib import cm
from imageio import imread, imwrite

imagename = 'balloons_resized_bw.jpg'
changed_pixels=40000
# Read the image file imagename into an array of numbers, IM
# Multiply by 1. / 255 to change the values so that they are floating point
# numbers ranging from 0 to 1.
IM = imread(imagename, as_gray=True) * (1. / 255)
IM_x, IM_y = IM.shape

for lost in range(changed_pixels):
```

```

x_,y_ = randint(1,IM_x-2), randint(1,IM_y-2)
val = .1*randn() + .5
IM[x_,y_] = max( min(val,1.), 0.)
imwrite("noised_" + imagename, IM)

```

A color image can be represented by three functions  $u_1, u_2$ , and  $u_3$ . In this lab we will work with black and white images, but total variation techniques can easily be used on more general images.

## A simple approach to image processing

Here is a first attempt at denoising: given a noisy image  $f$ , we look for a denoised image  $u$  minimizing the energy functional

$$J[u] = \int_{\Omega} L(x, u, \nabla u) dx, \quad (18.5)$$

where

$$\begin{aligned} L(x, u, \nabla u) &= \frac{1}{2}(u - f)^2 + \frac{\lambda}{2}|\nabla u|^2, \\ &= \frac{1}{2}(u - f)^2 + \frac{\lambda}{2}(u_x^2 + u_y^2)^2. \end{aligned}$$

This energy functional penalizes 1) images that are too different from the original noisy image, and 2) images that have large derivatives. The minimizing denoised image  $u$  will balance these two different costs.

Solving for the original denoised image  $u$  is a difficult inverse problem-some information is irretrievably lost when noise is introduced. However, a priori information can be used to guess at the structure of the original image. For example, here  $\lambda$  represents our best guess on how much noise was added to the image, and is known as a regularization parameter in inverse problem theory.

The Euler-Lagrange equation corresponding to (18.5) is

$$\begin{aligned} L_u - \operatorname{div} L_{\nabla u} &= (u - f) - \lambda \Delta u, \\ &= 0. \end{aligned}$$

and the gradient descent flow is

$$\begin{aligned} u_t &= -(u - f - \lambda \Delta u), \\ u(x, 0) &= f(x). \end{aligned} \quad (18.6)$$

Let  $u_{ij}^n$  represent our approximation to  $u(x_i, y_j)$  at time  $t_n$ . We will approximate  $u_t$  with a forward Euler difference, and  $\Delta u$  with centered differences:

$$\begin{aligned} u_t &\approx \frac{u_{ij}^{n+1} - u_{ij}^n}{\Delta t}, \\ u_{xx} &\approx \frac{u_{i+1,j}^n - 2u_{ij}^n + u_{i-1,j}^n}{\Delta x^2}, \\ u_{yy} &\approx \frac{u_{i,j+1}^n - 2u_{ij}^n + u_{i,j-1}^n}{\Delta y^2}. \end{aligned}$$



Original image

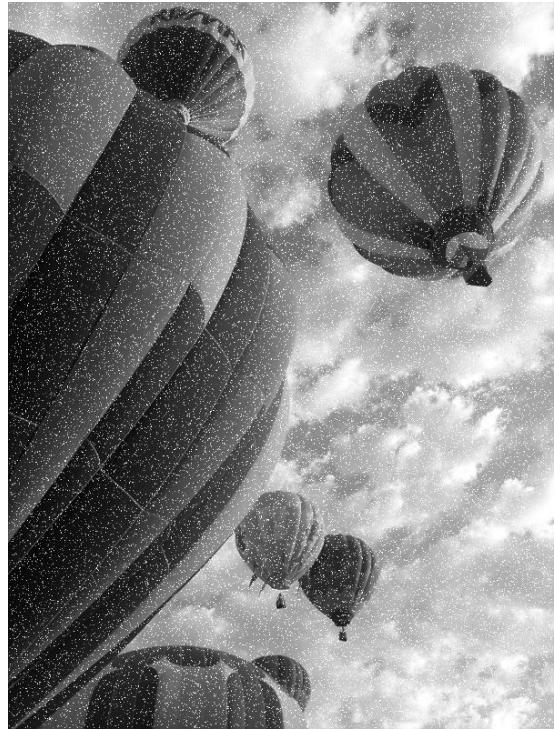


Image with white noise

Figure 18.2: Noise.

**Problem 2.** Using  $\Delta t = 1e-3$ ,  $\lambda = 40$ ,  $\Delta x = 1$ , and  $\Delta y = 1$ , implement the numerical scheme mentioned above to obtain a solution  $u$ . (So  $\Omega = [0, n_x] \times [0, n_y]$ , where  $n_x$  and  $n_y$  represent the number of pixels in the  $x$  and  $y$  dimensions, respectively.) Take 250 steps in time. Plot the original image as well as the image with noise. Compare your results with Figure 18.3.

Hint: Use the function `np.roll` to compute the spatial derivatives. For example, the second derivative can be approximated at interior grid points using

```
u_xx = np.roll(u,-1,axis=1) - 2*u + np.roll(u,1,axis=1)
```

## Image Processing: Total Variation Method

We represent an image by a function  $u : [0, 1] \times [0, 1] \rightarrow \mathbb{R}$ . A  $C^1$  function  $u : \Omega \rightarrow \mathbb{R}$  has bounded total variation on  $\Omega$  ( $BV(\Omega)$ ) if  $\int_{\Omega} |\nabla u| < \infty$ ;  $u$  is said to have total variation  $\int_{\Omega} |\nabla u|$ . Intuitively, the total variation of an image  $u$  increases when noise is added.

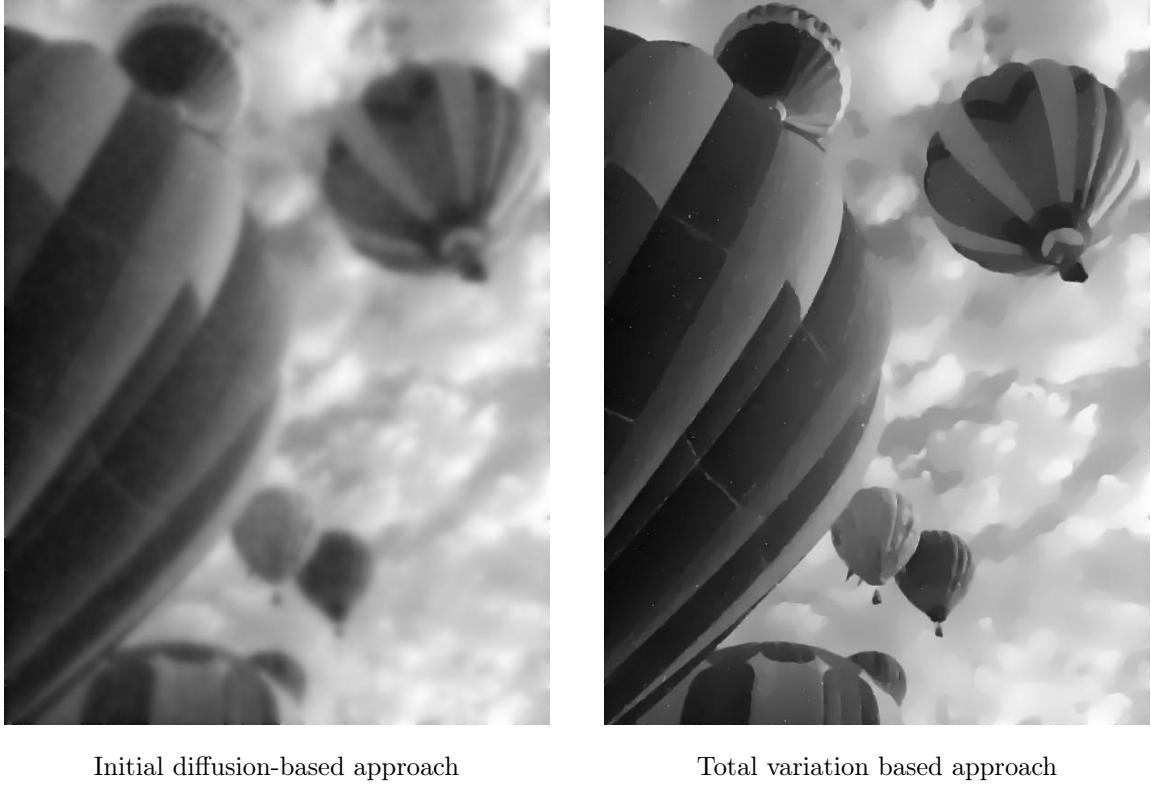


Figure 18.3: The solutions of (18.6) and (18.11), found using a first order Euler step in time and centered differences in space.

The total variation approach was originally introduced by Rudin, Osher, and Fatemi<sup>1</sup>. It was formulated as follows: given a noisy image  $f$ , we look to find a denoised image  $u$  minimizing

$$\int_{\Omega} |\nabla u(x)| dx \quad (18.7)$$

subject to the constraints

$$\int_{\Omega} u(x) dx = \int_{\Omega} f(x) dx, \quad (18.8)$$

$$\int_{\Omega} |u(x) - f(x)|^2 dx = \sigma |\Omega|. \quad (18.9)$$

Intuitively, (18.7) penalizes fast variations in  $f$  - this functional together with the constraint (18.8) has a constant minimum of  $u = \frac{1}{|\Omega|} \int_{\Omega} u(x) dx$ . This is obviously not what we want, so we add a constraint (18.9) specifying how far  $u(x)$  is required to differ from the noisy image  $f$ . More precisely, (18.8) specifies that the noise in the image has zero mean, and (18.9) requires that a variable  $\sigma$  be chosen a priori to represent the standard deviation of the noise.

Chambolle and Lions proved that the model introduced by Rudin, Osher, and Fatemi can be formulated equivalently as

$$F[u] = \min_{u \in BV(\Omega)} \int_{\Omega} |\nabla u| + \frac{\lambda}{2} (u - f)^2 dx, \quad (18.10)$$

---

<sup>1</sup>L. Rudin, S. Osher, and E. Fatemi, “Nonlinear total variation based noise removal algorithms”, *Physica D.*, 1992.

where  $\lambda > 0$  is a fixed regularization parameter<sup>2</sup>. Notice how this functional differs from (18.5):  $\int_{\Omega} |\nabla u|$  instead of  $\int_{\Omega} |\nabla u|^2$ . This turns out to cause a huge difference in the result. Mathematically, there is a nice way to extend  $F$  and the class of functions with bounded total variation to functions that are discontinuous across hyperplanes. The term  $\int |\nabla|$  tends to preserve edges/boundaries of objects in an image.

The gradient descent flow is given by

$$u_t = -\lambda(u - f) + \frac{u_{xx}u_y^2 + u_{yy}u_x^2 - 2u_xu_yu_{xy}}{(u_x^2 + u_y^2)^{3/2}}, \quad (18.11)$$

$$u(x, 0) = f(x).$$

Notice the singularity that occurs in the flow when  $|\nabla u| = 0$ . Numerically we will replace  $|\nabla u|^3$  in the denominator with  $(\varepsilon + |\nabla u|^2)^{3/2}$ , to remove the singularity.

**Problem 3.** Using  $\Delta t = 1e-3$ ,  $\lambda = 1$ ,  $\Delta x = 1$ , and  $\Delta y = 1$ , implement the numerical scheme mentioned above to obtain a solution  $u$ . Take 200 steps in time. Display both the diffusion-based and total variation images of the balloon. Compare your results with Figure 18.3. How small should  $\varepsilon$  be?

Hint: To compute the spatial derivatives, consider the following:

```
u_x = (np.roll(u,-1,axis=1) - np.roll(u,1,axis=1))/2
u_xx = np.roll(u,-1,axis=1) - 2*u + np.roll(u,1,axis=1)
u_xy = (np.roll(u_x,-1,axis=0) - np.roll(u_x,1,axis=0))/2.
```

---

<sup>2</sup>A. Chambolle and P.-L. Lions, "Image recovery via total variation minimization and related problems", Numer. Math., 1997.

# 19

## Transit time crossing a river

**Lab Objective:** This lab discusses a classical calculus of variations problem: how is a river to be crossed in the shortest possible time? We will look at a numerical solution using the pseudospectral method.

Suppose a boat is to be rowed across a river, from a point  $A$  on one side of a river ( $x = -1$ ), to a point  $B$  on the other side ( $x = 1$ ). Assuming the boat moves at a constant speed 1 relative to the current, how must the boat be steered to minimize the time required to cross the river?

Let us consider a typical trajectory for the boat as it crosses the river. If  $T$  is the time required to cross the river, then the position  $s$  of the boat at time  $t$  is

$$\begin{aligned} s(t) &= \langle x(t), y(t) \rangle, \quad t \in [0, T], \\ s'(t) &= \langle x'(t), y'(t) \rangle, \\ &= \langle \cos \theta(x(t)), \sin \theta(x(t)) \rangle + \langle 0, c(x(t)) \rangle. \end{aligned}$$

Here  $\langle \cos \theta, \sin \theta \rangle$  represents the motion of the boat due to the rower, and  $\langle 0, c \rangle$  is the motion of the boat due to the current.

We can relate the angle at which the boat is steered to the graph of its trajectory by noting that

$$\begin{aligned} y'(x) &= \frac{y'(t)}{x'(t)}, \\ &= \frac{\sin \theta + c}{\cos \theta}, \\ &= c \sec \theta + \tan \theta. \end{aligned} \tag{19.1}$$

The time  $T$  required to cross the river is given by

$$\begin{aligned} T &= \int_{-1}^1 t'(x) dx, \\ &= \int_{-1}^1 \frac{1}{x'(t)} dx \\ &= \int_{-1}^1 \sec(\theta) dx. \end{aligned} \tag{19.2}$$

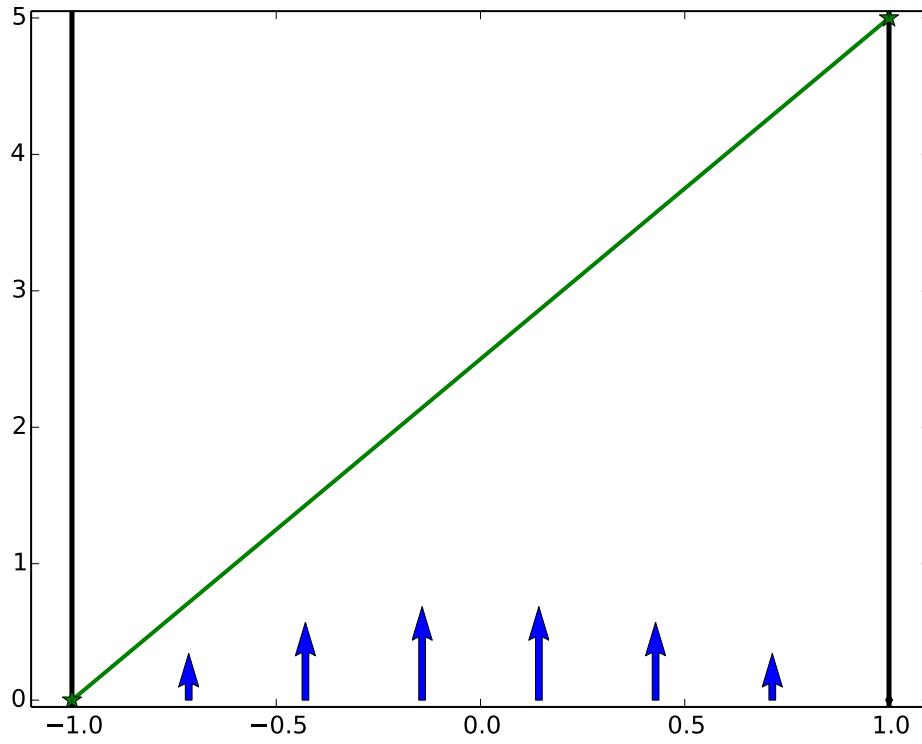


Figure 19.1: The river's current, along with a possible trajectory for the boat.

We would like to find an expression for the total time  $T$  required to cross the river from  $A$  to  $B$ , in terms of the graph of the boat's trajectory. To derive the functional  $T[y]$ , we note that

$$\begin{aligned} T[y] &= \int_{-1}^1 \sec \theta \, dx, \\ &= \int_{-1}^1 \frac{1}{1 - c^2} (c \tan \theta + \sec \theta - c^2 \sec \theta - c \tan \theta) \, dx, \\ &= \int_{-1}^1 \frac{1}{1 - c^2} (c \tan \theta + \sec \theta - cy') \, dx. \end{aligned}$$

Since

$$\begin{aligned} c \tan \theta + \sec \theta &= \sqrt{1 - c^2 + (c \sec \theta + \tan \theta)^2}, \\ &= \sqrt{1 - c^2 + (y')^2}, \end{aligned}$$

we obtain at last

$$T[y] = \int_{-1}^1 \left[ \alpha(x) \sqrt{1 + (\alpha y')^2(x)} - (\alpha^2 c y')(x) \right] dx, \quad (19.3)$$

where  $\alpha = (1 - c^2)^{-1/2}$ .

**Problem 1.** Assume that the current is given by  $c(x) = -\frac{7}{10}(x^2 - 1)$ . (This function assumes, for example, that the current is faster near the center of the river.) Write 2 python functions. The first should accept as arguments a function  $y$ , its derivative  $y'$ , and an  $x$ -value, and return  $L(x, y(x), y'(x))$  (where  $T[y] = \int_{-1}^1 L(x, y(x), y'(x)) dx$ ). The second should use the first function to compute and return  $T[y]$  for a given path  $y(x)$ .

(Hint: The integration for  $T[y]$  can be done use an approximation method such as the midpoint method or can be done using the `quad` function from `scipy.integrate`.)

**Problem 2.** Let  $y(x)$  be the straight-line path between  $A = (-1, 0)$  and  $B = (1, 5)$ . Numerically calculate  $T[y]$  to get an upper bound on the minimum time required to cross from  $A$  to  $B$ . Using (19.2), find a lower bound on the minimum time required to cross.

(Hint: if  $G = \int f(x)dx$  and we want to minimize  $G$ , try minimizing  $f(x)$ .)

We look for the path  $y(x)$  that minimizes the time required for the boat to cross the river, so that the function  $T$  is minimized. From the calculus of variations we know that a smooth path  $y(x)$  minimizes  $T$  only if the Euler-Lagrange equation is satisfied. Recall that the Euler-Lagrange equation is

$$L_y - \frac{d}{dx} L_{y'} = 0.$$

Since  $L_y = 0$ , we see that the shortest time trajectory satisfies

$$\frac{d}{dx} L_{y'} = \frac{d}{dx} \left( \alpha^3(x)y'(x)(1 + (\alpha y')^2(x))^{-1/2} - \alpha^2(x)c \right) = 0. \quad (19.4)$$

**Problem 3.** Numerically solve the Euler-Lagrange equation (19.4), using  $c(x) = -\frac{7}{10}(x^2 - 1)$  and  $\alpha = (1 - c^2)^{-1/2}$ , and  $y(-1) = 0$ ,  $y(1) = 5$ .

Hint: Since this boundary value problem is defined over the domain  $[-1, 1]$ , it is easy to solve using the pseudospectral method. Begin by replacing each  $\frac{d}{dx}$  with the pseudospectral differentiation matrix  $D$ . Then impose the boundary conditions and solve.

**Problem 4.** Plot the angle at which the boat should be pointed at each  $x$ -coordinate. (Hint: Use Equation (19.1); see Figure 19.3. Note that the angle the boat should be steered is not described by the tangent vector to the trajectory. Consider using `scipy.optimize.root` or `scipy.interpolate.barycentric_interpolate`)

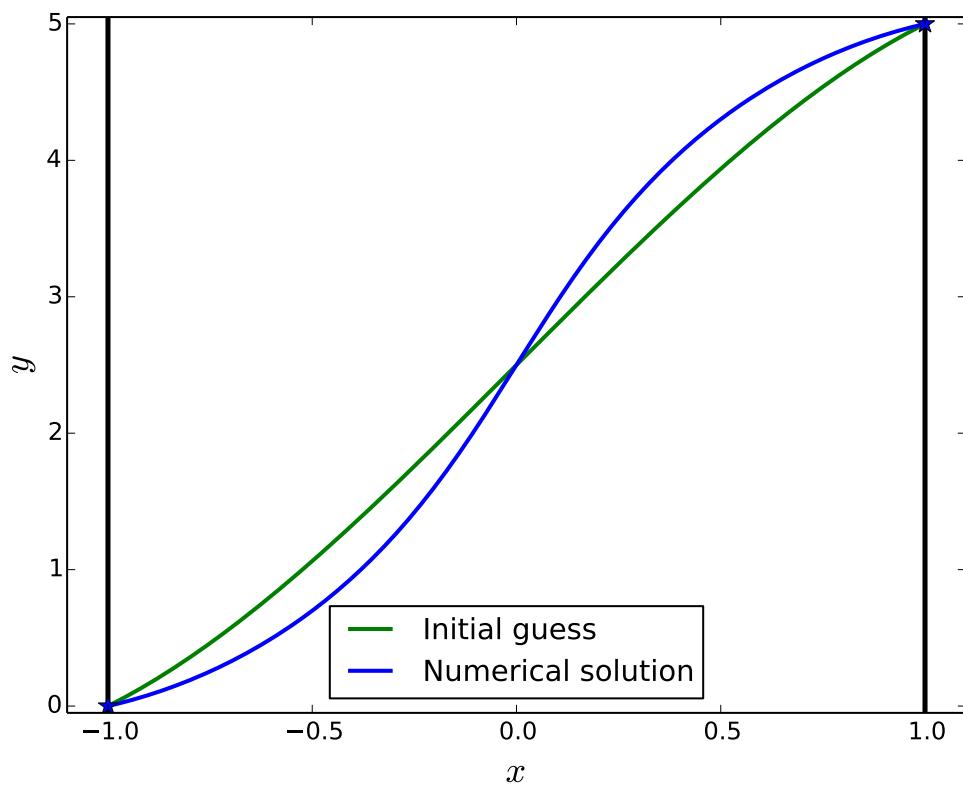


Figure 19.2: Numerical computation of the trajectory with the shortest transit time.

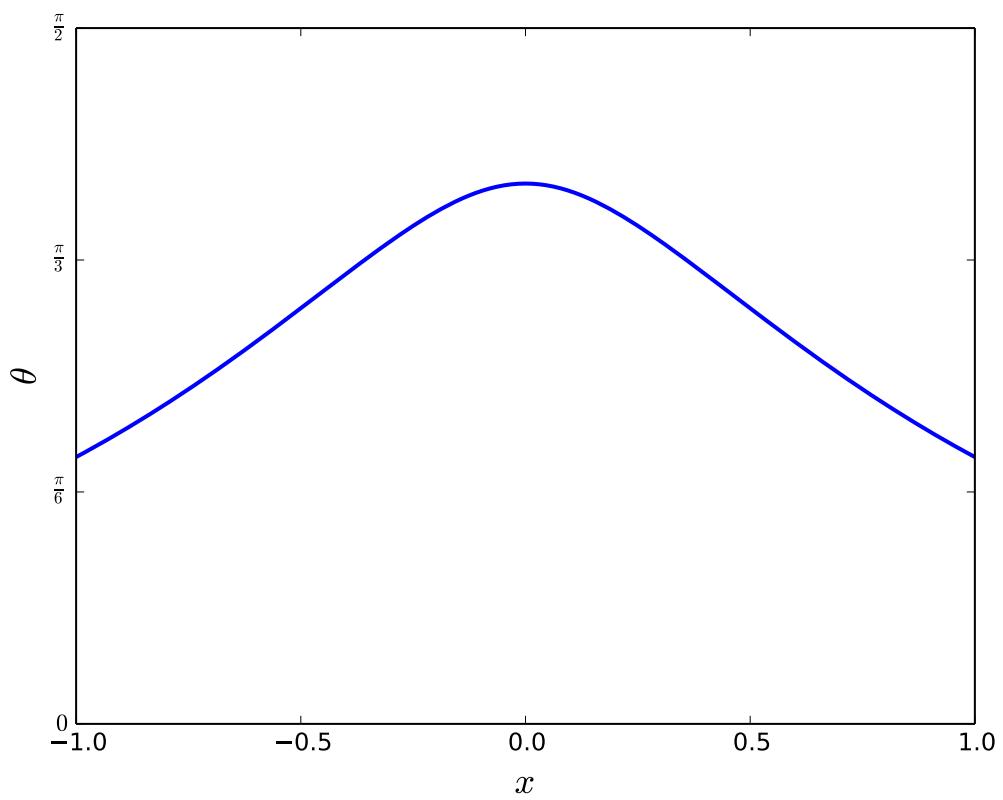


Figure 19.3: The optimal angle to steer the boat.



# 20 HIV Treatment Using Optimal Control

## Introduction

Viruses cause many common illnesses in society today, including influenza, the common cold, and COVID-19. Viruses are obligate parasites, meaning that they must infect a host in order to replicate. After entering a host cell, viruses hijack host machinery to replicate their genome and translate their proteins. After this process, the new virus particles are assembled and lyse (break apart) the host cell to find a new host.

Mammalian immune systems are composed of two interconnected systems: the innate immune system and the adaptive immune system. While both branches of the immune system can combat viruses, the adaptive immune system is especially suited to recognize and neutralize viral infections. A major part of the adaptive immune response is helper T cells, as these cells moderate and regulate all other facets of the immune response. Helper T cells are most characterized by the presence of a receptor called CD4, which helps the cell recognize infections.

One of the most devastating viral illnesses today is acquired immunodeficiency syndrome (AIDS), caused by the human immunodeficiency virus (HIV). HIV specifically targets and replicates in helper T cells, rendering them nonfunctional and killing them. By taking out the most important regulator of the immune system, HIV makes it difficult for the body to fight infection, so sicknesses that would normally be trivial for the body to manage, such as the common cold, yeast infections, and pneumonia, become deadly.

Currently, there is no cure for HIV, and vaccines are difficult to develop. Treatments that curb the replication of HIV and help maintain healthy helper T cell population levels are available, but they are expensive and must be taken for the rest of a patient's life. Optimizing the dosage is essential to maximize the drug's effect while minimizing the cost and negative side-effects of long-term usage. In this lab, we will use optimal control to find the optimum dosage of a two-drug combination to fight HIV. In this lab we will use optimal control to find the optimal dosage of a two-drug combination<sup>1</sup>.

---

<sup>1</sup>SHORT COURSES ON THE MATHEMATICS OF BIOLOGICAL COMPLEXITY, Web. 15 Apr. 2015  
<http://www.math.utk.edu/~lenhart/smb2003.v2.html>.

## Derivation of Control

We begin by defining some variables. Let  $T$  represents the concentration of  $CD4^+T$  cells and  $V$  the concentration of HIV particles.  $s_1$  and  $s_2$  represent the production of T cells by various processes.  $B_1$  and  $B_2$  are half saturation constants (sort of like crowd control in the blood stream and plasma). Let  $\mu$  be the death rate of uninfected T cells,  $k$  the rate of infection of T cells, and  $c$  the death rate of the virus. Let  $g$  be the input rate of some external viral source. The control variables  $u_1$  and  $u_2$  represent the amount of drugs that introduce new T cells or kill the virus, respectively.<sup>2</sup>

Next we write the state system, the equations that describe the changes in T cells and viruses:

$$\begin{aligned}\frac{dT(t)}{dt} &= s_1 - \frac{s_2 V(t)}{B_1 + V(t)} - \mu T(t) - kV(t)T(t) + u_1(t)T(t), \\ \frac{dV(t)}{dt} &= \frac{gV(t)}{B_2 + V(t)}(1 - u_2(t)) - cV(t)T(t).\end{aligned}\tag{20.1}$$

The term  $s_1 - \frac{s_2 V(t)}{B_1 + V(t)}$  is the source/proliferation of unaffected T cells,  $\mu T(t)$  the natural loss of T cells,  $kV(t)T(t)$  the loss of T cells by infection.  $\frac{gV(t)}{B_2 + V(t)}$  represents the viral contribution to plasma, and  $cV(t)T(t)$  the viral loss. To these equations we add initial conditions  $T(0) = T_0$  and  $V(0) = V_0$ .<sup>3</sup>

We now seek to maximize the functional

$$J(u_1, u_2) = \int_0^{t_f} [T - (A_1 u_1^2 + A_2 u_2^2)] dt.$$

This functional considers i) the benefit of T cells, and ii) the systematic costs of drug treatments. The constants  $A_1$  and  $A_2$  represent scalars to adjust the size of terms coming from  $u_1^2$  and  $u_2^2$  respectively. We seek an optimal control  $u_1^*, u_2^*$  satisfying

$$J(u_1^*, u_2^*) = \max_{(u_1, u_2) \in U} J(u_1, u_2) = \min_{(u_1, u_2) \in U} -J(u_1, u_2),$$

where  $U = \{(u_1, u_2) : a_i \leq u_i(t) \leq b_i \text{ for } t \in [0, t_f], i = 1, 2\}$ .

## Optimality System

The Hamiltonian is defined as:

$$\begin{aligned}H &= \vec{\lambda} \cdot \vec{f} - L \\ H &= \lambda_1 \left[ s_1 - \frac{s_2 V}{B_1 + V} - \mu T - kVT + u_1 T \right] + \lambda_2 \left[ \frac{g(1 - u_2)V}{B_2 + V} - cVT \right] \\ &\quad + [T - (A_1 u_1^2 + A_2 u_2^2)].\end{aligned}$$

Note that the costate is represented with  $\lambda$  instead of  $p$ . The costate evolution equations are:

$$\begin{aligned}\lambda_1' &= -\frac{\partial H}{\partial T} = -1 + \lambda_1[\mu + kV^* - u_1^*] + \lambda_2 cV^*, \\ \lambda_2' &= -\frac{\partial H}{\partial V} = \lambda_1 \left[ \frac{B_1 s_2}{(B_1 + V^*)^2} + kT^* \right] - \lambda_2 \left[ \frac{B_2 g(1 - u_2^*)}{(B_2 + V^*)^2} - cT^* \right]\end{aligned}$$

<sup>2</sup>'Immunotherapy of HIV-1 Infection', Kirschner, D. and Webb, G. F., Journal of Biological Systems, 6(1), 71-83 (1998)

<sup>3</sup>'Optimal Control of an HIV Immunology Model', H.R Joshi

where  $T^*, V^*$  denote the optimal  $T, V$ . The endpoint conditions are  $\lambda_1(t_f) = \lambda_2(t_f) = 0$ , with  $T(0) = T_0$  and  $V(0) = V_0$ . Using Pontryagin's maximum principle to find the control, we have

$$\frac{\partial H}{\partial u_1} = -2A_1 u_1^*(t) + \lambda_1 T^*(t) = 0 \quad \frac{\partial H}{\partial u_2} = -2A_2 u_2^*(t) + \lambda_2 \left[ \frac{-gV^*(t)}{B_2 + V^*(t)} \right] = 0$$

which gives (provided these are within the bounds of the controls)

$$u_1^*(t) = \frac{1}{2A_1} [\lambda_1 T^*(t)], \\ u_2^*(t) = \frac{-1}{2A_2} \left[ \lambda_2 \frac{gV^*(t)}{B_2 + V^*(t)} \right].$$

Taking into account the bounds on the controls, we have

$$u_1^*(t) = \min \left\{ \max \left\{ a_1, \frac{1}{2A_1} (\lambda_1 T^*(t)) \right\}, b_1 \right\}, \\ u_2^*(t) = \min \left\{ \max \left\{ a_2, \frac{-\lambda_2}{2A_2} \frac{gV^*(t)}{B_2 + V^*(t)} \right\}, b_2 \right\}.$$

This gives us the optimal system

$$T' = s_1 - \frac{s_2 V}{B_1 + V} - \mu T - k V T + \min \left\{ \max \left\{ a_1, \frac{1}{2A_1} (\lambda_1 T) \right\}, b_1 \right\} T, \\ V' = \frac{g(1 - \min \left\{ \max \left\{ a_2, \frac{-\lambda_2}{2A_2} \frac{gV}{B_2 + V} \right\}, b_2 \right\}) V}{B_2 + V} - c V T \quad (20.2)$$

$$\lambda'_1 = -1 + \lambda_1 \left[ \mu + k V - \min \left\{ \max \left\{ a_1, \frac{1}{2A_1} (\lambda_1 T) \right\}, b_1 \right\} \right] + \lambda_2 c V, \\ \lambda'_2 = \lambda_1 \left[ \frac{B_1 s_2}{(B_1 + V)^2} + k T \right] - \lambda_2 \left[ \frac{B_2 g(1 - \min \left\{ \max \left\{ a_2, \frac{-\lambda_2}{2A_2} \frac{V}{B_2 + V} \right\}, b_2 \right\})}{(B_2 + V)^2} - c T \right], \quad (20.3)$$

with end conditions  $\lambda_1(t_f) = \lambda_2(t_f) = 0$ , and  $T(0) = T_0, V(0) = V_0$ .

## Creating a Numerical Solver

We iteratively solve for our control  $u$ . In each iteration we solve our state equations and our costate equations numerically, then use those to find our new control. Lastly, we check to see if our control has converged. To solve each set of differential equations, we will use the RK4 solver from a previous lab with one minor adjustment. Our state equations depend on  $u$ , and our costate equations depend on our state equations. Therefore, we will pass another parameter into the function that RK4 takes in that will index the arrays our equations depend on.

```
# Dependencies for this lab's code:
import numpy as np
from matplotlib import pyplot as plt

#Code from RK4 Lab with minor edits
def initialize_all(y0, t0, tf, n):
    """ An initialization routine for the different ODE solving
    methods in the lab. This initializes Y, T, and h. """

```

```

if isinstance(y0, np.ndarray):
    Y = np.empty((n, y0.size)).squeeze()
else:
    Y = np.empty(n)
Y[0] = y0
T = np.linspace(t0, tf, n)
h = float(tf - t0) / (n - 1)
return Y, T, h

def RK4(f, y0, t0, tf, n):
    """ Use the RK4 method to compute an approximate solution
    to the ODE  $y' = f(t, y)$  at n equispaced parameter values from t0 to t
    with initial conditions  $y(t0) = y0$ .

    y0 is assumed to be either a constant or a one-dimensional numpy array.
    tf and t0 are assumed to be constants.
    f is assumed to accept three arguments.
    The first is a constant giving the value of t.
    The second is a one-dimensional numpy array of the same size as y.
    The third is an index to the other arrays.

    This function returns an array Y of shape (n,) if
    y is a constant or an array of size 1.
    It returns an array of shape (n, y.size) otherwise.
    In either case, Y[i] is the approximate value of y at
    the i'th value of np.linspace(t0, tf, n).
    """
    Y,T,h = initialize_all(y0,t0,tf,n)
    for i in range(n-1):
        K1 = f(T[i],Y[i],i)
        K2 = f(T[i]+h/2.,Y[i]+h/2.*K1,i)
        K3 = f(T[i]+h/2.,Y[i]+h/2.*K2,i)
        K4 = f(T[i+1],Y[i]+h*K3,i)
        Y[i+1] = Y[i] + h/6.*(K1+2*K2 +2*K3+K4)
    return Y

```

**Problem 1.** Create a function that defines the state equations and returns both equations in a single array. The function should be able to be passed into the RK4 solver. This function can depend on the global variables defined below.

### Achtung!

When solving the state equations, because of the nature of  $T'$  and  $V'$ , solve the original equations (20.1) from the beginning of the lab and not the equations (20.2) with  $u_i^*(t)$  replaced by the minmax function.

```
a_1, a_2 = 0, 0
b_1, b_2 = 0.02, 0.9
s_1, s_2 = 2., 1.5
mu = 0.002
k = 0.000025
g = 30.
c = 0.007
B_1, B_2 = 14, 1
A_1, A_2 = 250000, 75
T0, V0 = 400, 3
t_f = 50
n = 2000
```

These constants come from both references cited at the end of this lab.

```
# initialize global variables, state, costate, and u.
state = np.zeros((n,2))
state0 = np.array([T0, V0])

costate = np.zeros((n,2))
costate0 = np.zeros(2)

u=np.zeros((n,2))
u[:,0] = .02
u[:,1] = .9

# define state equations
def state_equations(t,y,i):
    ...
    Parameters
    -----
    t : float
        the time
    y : ndarray (2,)
        the T cell concentration and the Virus concentration at time t
    i : int
        index for the global variable u.
    Returns
    -----
    y_dot : ndarray (2,)
```

```

the derivative of the T cell concentration and the virus ←
concentration at time t
...
pass
```

The state equations work great in the RK4 solver; however, the costate equations have end conditions rather than initial conditions. Thus we want our RK4 solver to iterate backwards from the end to the beginning. An easy way to accomplish this is to define a function  $\hat{\lambda}_i(t) = \lambda_i(t_f - t)$ . Then  $\hat{\lambda}_i$  has the initial conditions  $\hat{\lambda}_i(0) = \lambda_i(t_f)$ . We get the new equations

$$\begin{aligned}\hat{\lambda}'_1(t) &= \lambda_1(t_f - t) (-\mu - kV(t_f - t) + u_1(t_f - t)) - c\lambda_2(t_f - t)V(t_f - t) + 1, \\ \hat{\lambda}'_2(t) &= -\lambda_1(t_f - t) \left( \frac{s_2 B_1}{(B_1 + V(t_f - t))^2} + kT(t_f - t) \right) \\ &\quad + \lambda_2(t_f - t) \left( \frac{g B_2 (1 - u_2(t_f - t))}{(B_2 + V(t_f - t))^2} - cT(t_f - t) \right).\end{aligned}$$

These we can solve with our RK4 solver and recover the original costate equations by simply indexing the array backwards.

**Problem 2.** Create a function that defines the costate equations and returns both equations in a single array. The function should be able to be passed into the RK4 solver. Use the global variables as defined in Problem 1.

```

def lambda_hat(t,y,i):
    ...
    Parameters
    -----
    t : float
        the time
    y : ndarray (2,)
        the lambda_hat values at time t
    i : int
        index for global variables, u and state.
    Returns
    -----
    y_dot : ndarray (2,)
        the derivative of the lambda_hats at time t.
    ...
    pass
```

Finally, we can put these together to create our solver.

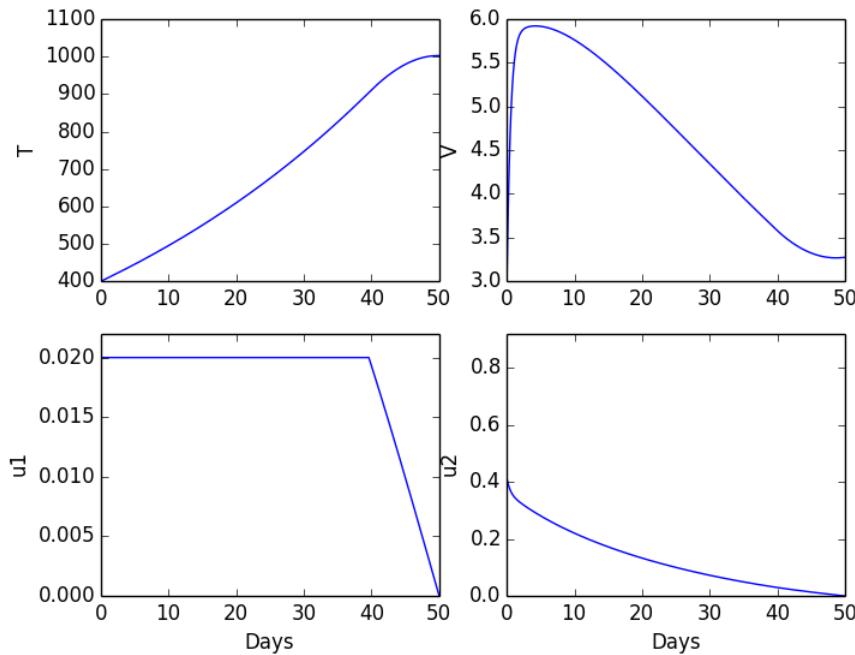


Figure 20.1: The solution to Problem 3.

**Problem 3.** Create and run a numerical solver for the HIV two drug model using the code below.

```

epsilon = 0.001
test = epsilon + 1

while(test > epsilon):
    oldu = u.copy();

    #solve the state equations with forward iteration
    #state = RK4(...)

    #solve the costate equations with backwards iteration
    #costate = RK4(...)[::-1]

    #solve for u1 and u2

    #update control
    u[:,0] = 0.5*(u1 + oldu[:,0])
    u[:,1] = 0.5*(u2 + oldu[:,1])

    #test for convergence

```

```
test = abs(oldu - u).sum()
```

Your solutions should match Figure 20.1.

Patients usually take several different classes of drugs at a time to prevent HIV from replicating and progressing into AIDS. Reverse transcriptase inhibitors prevent the HIV genome from inserting itself into the host genome. These prevent helper T cell death by lowering the number of HIV particles in the body. Protease inhibitors prevent the activation of HIV proteins that are needed for replication. Fusion inhibitors can be taken early in the course of HIV infection and prevent the entry of HIV into helper T cells. There are many unique drugs in each class, all with known and unknown interactions and side effects. Physicians rotate through drugs to help their patients have a positive outcome and to prevent the virus from becoming resistant to any one drug.

# 21 Solitons

**Lab Objective:** Use a pseudospectral method to study solitons, the traveling wave solutions of the Korteweg-de Vries equation.

The Korteweg-de Vries (KdV) equation is a partial differential equation given by

$$\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + \frac{\partial^3 u}{\partial x^3} = 0.$$

that describes shallow water waves.

The KdV equation possesses traveling wave solutions called solitons. These traveling waves have the form

$$u(x, t) = 3s \operatorname{sech}^2 \left( \frac{\sqrt{s}}{2}(x - st - a) \right),$$

where  $s$  is the speed of the wave. Solitons were first studied by John Scott Russell in 1834, in the Union Canal in Scotland. When a canal boat suddenly stopped, the water piled up in front of the boat continued moving down the canal in the shape of a pulse.

Note that there is a soliton solution for each wave speed  $s$ , and that the amplitude and speed of the soliton determine each other. Solitons are nonlinearly stable (bumped waves return to their previous shape), and they maintain their energy as they travel. Two interacting solitons will also both maintain their shapes after crossing paths.

## Numerical solution

Consider the KdV equation on  $[-\pi, \pi]$ , together with an appropriate initial condition:

$$\begin{aligned} u_t &= -\frac{1}{2} (u^2)_x - u_{xxx}, \\ u(x, 0) &= u_0(x). \end{aligned}$$

This form of the equation is slightly more convenient for the approach we will take. We will suppose the initial condition is equal at the two endpoints; that is,  $u_0(-\pi) = u_0(\pi)$ . This will allow us to use the pseudospectral method to find a numerical approximation for the solution  $u(x, t)$ .

As a reminder, the pseudospectral method involves writing the solution at each point in time using a set of basis functions, complex exponentials being the most common, and using this representation to convert the PDE into an ODE. Specifically, we can write any solution  $u(x, t)$  as

$$u(x, t) = \sum_{k=-\infty}^{\infty} y_k(t) e^{ikx}.$$

Recall that  $k$  is known as the wave number. Note that all time-dependence of the solution is contained in the coefficients. We can only compute this to some finite precision, so we will choose some  $n$  and truncate the series as

$$u(x, t) = \sum_{k=-n}^n y_k(t) e^{ikx}.$$

The objective is to obtain an ordinary differential equation for the coefficients  $y_k(t)$ . We now plug it into the PDE:

$$\begin{aligned} \frac{\partial}{\partial t} \sum_{k=-n}^n y_k(t) e^{ikx} &= -\frac{\partial}{\partial x} \left( \sum_{k=-n}^n y_k(t) e^{ikx} \right)^2 - \frac{\partial^3}{\partial^3 x} \sum_{k=-n}^n y_k(t) e^{ikx} \\ \sum_{k=-n}^n y'_k(t) e^{ikx} &= -\frac{\partial}{\partial x} \left( \sum_{k=-n}^n y_k(t) e^{ikx} \right)^2 + \sum_{k=-n}^n ik^3 y_k(t) e^{ikx} \end{aligned}$$

For this particular PDE, this leads to an apparent problem: the  $u^2$  term will be difficult and computationally costly to differentiate. However, we can get around this difficulty using the fast Fourier transform.

Divide  $[-\pi, \pi]$  into  $2n+1$  intervals of equal width  $\frac{2\pi}{2n+1}$ , and let  $-\pi = x_{-n}, x_{-n+1}, \dots, x_n, x_{n+1} = \pi$  be the  $2n+2$  evenly-spaced gridpoints. For any function  $f$  on that interval with Fourier series  $f(x) = \sum_{k=-\infty}^{\infty} a_k e^{ikx}$ , we can use the discrete Fourier transform on the values  $f(x_{-n}), \dots, f(x_{n+1})$  at the gridpoints to quickly get the Fourier coefficients  $a_{-n}, \dots, a_n$ . The inverse Fourier transform can be used to get the function values at the grid points from the Fourier coefficients. Both of these operations are very efficient, having complexity  $O(n \log n)$ . This sets up our strategy.

At each time  $t$ , we can use the inverse Fourier transform to compute the values of  $u(x_m, t)$  for  $m = -n, \dots, n+1$ . Then, we apply the Fourier transform to  $u^2$  to get its Fourier coefficients. We will denote these as  $w_k$ , so

$$u^2(x, t) = \sum_{k=-n}^n w_k(t) e^{ikx}.$$

Then,

$$\frac{\partial}{\partial x} u^2(x, t) = \sum_{k=-n}^n ikw_k(t) e^{ikx},$$

so the KdV equation can be written as

$$\sum_{k=-n}^n y'_k(t) e^{ikx} = \sum_{k=-n}^n \left( -\frac{1}{2} ikw_k(t) + ik^3 y_k(t) \right) e^{ikx}$$

Equating terms in the Fourier series yields the ordinary system of differential equations

$$y'_k = -\frac{1}{2} ikw_k + ik^3 y_k, \quad k = -n, \dots, n.$$

We can also write this in a vectorized form as

$$\mathbf{y}' = -\frac{1}{2}i\mathbf{k}\mathcal{F}(\mathcal{F}^{-1}(\mathbf{y})^2) + ik^3\mathbf{y} \quad (21.1)$$

where  $\mathcal{F}$  denotes the discrete Fourier transform and multiplication of vectors is componentwise. To obtain the initial condition for the  $y_k$ , we can simply use the discrete Fourier transform again:

$$\mathbf{y}(0) = \mathcal{F}(u_0(x_{-n}), \dots, u_0(x_{n+1}))$$

To compute the fast Fourier and inverse fast Fourier transforms numerically, we will use the `scipy.fft` module, which has functions `fft` for the fast Fourier transform and `ifft` for the inverse fast Fourier transform. These functions use an order for the coefficients that is slightly nonintuitive: the coefficients for  $k \geq 0$  are all listed first, followed by the coefficients for  $k < 0$ . The vector of wavenumbers can be created as follows:

```
k = np.concatenate([
    np.arange(0,n+1),
    np.arange(-n-1,0)
])
```

We are now prepared to numerically solve the KdV equation.

**Problem 1.** Write a function that accepts the time value  $t$  (which won't be used here, but will be useful later) the vector  $\mathbf{y} = (y_0, y_1, \dots, y_n, y_{-n-1}, \dots, y_{-1})$  and the vector of  $k$  values and returns  $\mathbf{y}'$ .

To numerically solve this ODE, use the following implementation of the RK4 algorithm:

```
def RK4(f, y0, T, dt, k):
    """
    Solves the ODE y'=f(t, y) using the Runge-Kutta 4 method with initial
    condition y0 on the time interval [0,T] using a time step of dt.
    The value of k is passed directly into the function f.

    Returns:
        t ((T,) ndarray) - the time values
        Y ((T, 2n+2) ndarray) - the solution values. The solution at the
            i-th time step can be indexed as Y[i].
    """
    # Set up matrices for the solution
    ts = np.arange(0, T+dt, dt)
    Y = np.empty((len(ts), len(k)), dtype=complex)
    y = y0
    Y[0] = y
    for i in range(1, len(ts)):
        # Use RK4
        t = ts[i]
        K1 = f(t, y, k)
        K2 = f(t + dt/2, y + K1/2, k)
        K3 = f(t + dt/2, y + K2/2, k)
        K4 = f(t + dt, y + K3, k)
        y = y + (K1 + 2*K2 + 2*K3 + K4) * dt / 6
        Y[i] = y
```

```

K2 = f(t + dt/2, y + 0.5*dt*K1, k)
K3 = f(t + dt/2, y + 0.5*dt*K2, k)
K4 = f(t + dt, y + dt*K3, k)
y = y + (dt / 6.) * (K1 + 2*K2 + 2*K3 + K4)
Y[i] = y
return ts, Y

```

Once we have solved for the coefficients  $y(t)$ , we need to convert them back into function values  $u(x, t)$  in order to visualize the solution. This is accomplished by using the `ifft` function on the coefficient values at each time step. However, this function is designed to work with complex numbers, and returns a complex-valued array. Due to numerical error, even though our ODE solution is real-valued, there may be small imaginary components to the result; use `np.real` on the result to discard these.

**Problem 2.** Write a function that accepts an initial condition  $u_0$ , a final time  $T$ , the timestep  $dt$ , an integer  $n$  for the number of coefficients to use, and another integer `skip`. Numerically solve for the coefficients  $y(t)$  of a solution to the KdV equation.

Next, convert the Fourier coefficients back into function values at the gridpoints using the inverse Fourier transform. However, only do this for every `skip`-th timestep; we will be using far more timesteps than we want to plot. Return the gridpoints, the timesteps, and the solution  $u(x, t)$ .

Once we have the function values, we can plot them as a surface as follows:

```

fig = plt.figure()
ax = fig.add_subplot(1,1,1, projection='3d')

T, X = np.meshgrid(t, x, indexing='ij')
ax.plot_surface(T, X, u, cmap='coolwarm', rstride=1, cstride=1)

```

**Problem 3.** Numerically solve the KdV equation on  $[-\pi, \pi]$  with initial conditions

$$u(x, t = 0) = 3s \operatorname{sech}^2 \left( \frac{\sqrt{s}}{2}(x + a) \right),$$

where  $s = 25^2$ ,  $a = 2$ . Solve on the time domain  $[0, 0.0075]$ , and use  $n = 127$ . Compare with Figure 21.1; to get a similar viewpoint, use the following:

```

ax.view_init(elev=45, azim=-45)
ax.set_zlim(0, 4000)
ax.invert_xaxis()

```

How small of a timestep did you need to use for the numerical integration to be stable?

Hint: `numpy` does not have a `sech` function; use `1/cosh(x)` to compute it instead.

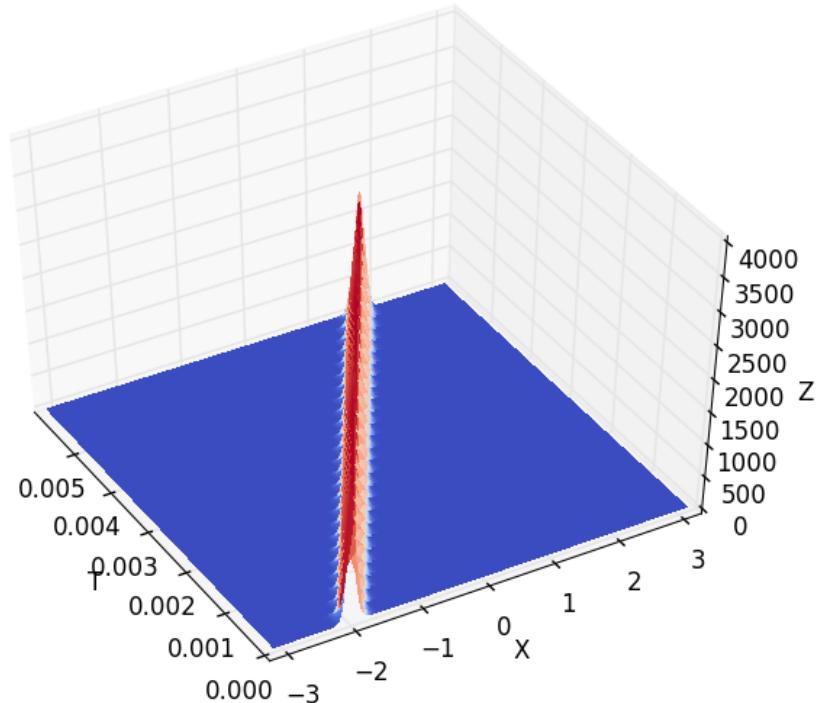


Figure 21.1: The solution to Problem 3.

**Problem 4.** Numerically solve the KdV equation on  $[-\pi, \pi]$ . This time we define the initial condition to be the superposition of two solitons:

$$u(x, t = 0) = 3s_1 \operatorname{sech}^2 \left( \frac{\sqrt{s_1}}{2}(x + a_1) \right) + 3s_2 \operatorname{sech}^2 \left( \frac{\sqrt{s_2}}{2}(x + a_2) \right),$$

where  $s_1 = 25^2$ ,  $a_1 = 2$ , and  $s_2 = 16^2$ ,  $a_2 = 1$ .<sup>a</sup> Solve on the time domain  $[0, 0.0075]$ . The solution is shown in Figure 21.2.

---

<sup>a</sup>This problem is from Spectral Methods in MATLAB, by Trefethen.

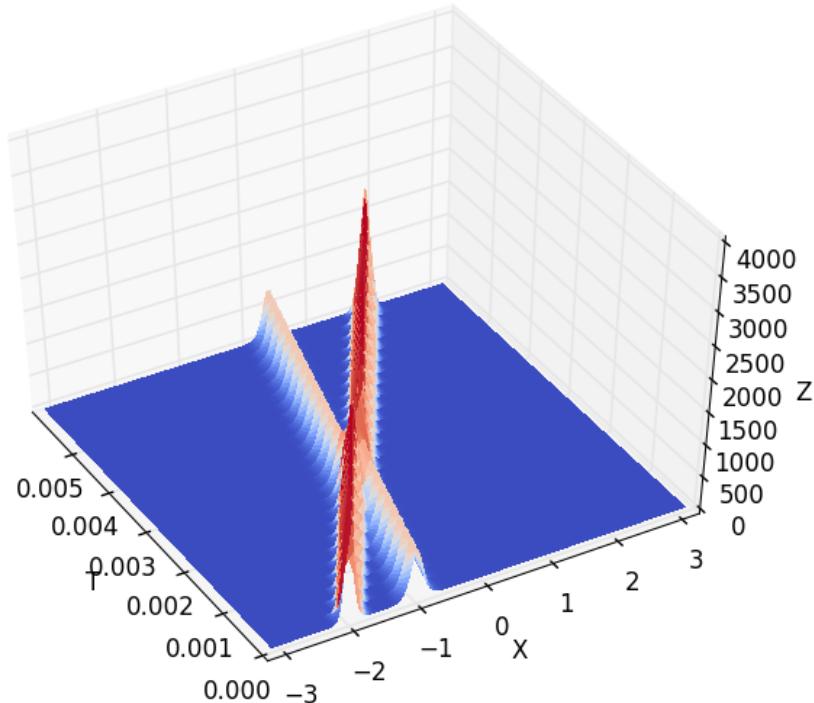


Figure 21.2: The solution to Problem 4.

**Problem 5.** Consider again equation (21.1). The linear term in this equation is  $ik^3\mathbf{y}$ . This term contributes much of the exponential growth in the ODE, and contributes to how short the time step must be to ensure numerical stability. Make the substitution  $z_k(t) = e^{-ik^3t}y_k(t)$  and find a similar ODE for  $\mathbf{z}$ . This essentially allows the exponential growth to be scaled out (it's solved for analytically, replacing it with rotation in the complex plane). Use the resulting equation to solve the previous problem. How much larger of a timestep can you use while this method remains stable?

# 22 Obstacle Avoidance

**Lab Objective:** Solve boundary value problems that arise when using Pontryagin's Maximum principle.

## Pontryagin's Maximum Principle

Now that we understand how to solve boundary value problems, we can apply this to solve optimal control problems. Pontryagin's Maximum Principle is a very common way to formulate control problems as BVPs.

### Fixed Time, Fixed Endpoint

We will begin with the more simple fixed time horizon problems. Fixed time horizon problems are commonly reformulated as boundary value problems, and we can apply what we have already learned about solving BVPs to make these problems easier to solve. We introduce fixed time horizon problems with a cost functional of the following form

$$J(u) = \int_{t_0}^{t_f} L(t, s(t), u(t)) dt + K(t_f, s_f), \quad (22.1)$$

where  $t_0$  and  $t_f$  are fixed. In this functional,  $L(t, s(t), u(t))$  represents the cost of a certain path determined by the control  $u$ , and  $K(t_f, s_f)$  is the terminal cost. We also have that

$$\dot{s} = f(t, s, u), \quad s_0 = s(t_0), \quad s_f = s(t_f). \quad (22.2)$$

In these equations  $t$  is time,  $s$  is the state variable, and  $u$  is the control variable. The maximum principle also uses the Hamiltonian equation

$$H(t, s, u, p) = \langle p, f(t, s, u) \rangle - L(t, s, u), \quad (22.3)$$

where  $p$  is a newly introduced variable called the costate. This Hamiltonian is then used to define an ODE system. This first equation defines a costate ODE system

$$\dot{p}^* = -H_s(t, s^*, u^*, p^*), \quad (22.4)$$

where a variable marked with an asterisk is the optimal choice of that variable, meaning that equation 22.4 is only true for the optimal state  $s^*$ , costate  $p^*$ , and control  $u^*$  functions. This next equation will allow us to solve for the control in terms of the state and costate

$$0 = H_u(t, s^*, u^*, p^*), \quad \forall t \in [t_0, t_f]. \quad (22.5)$$

The combination of these equations will allow us to create a BVP that will solve for the optimal control  $u^*$  and the associated states  $s^*$ . Our ODE comes from 22.2, 22.4, and 22.5, and the boundary values will come from our initial and final conditions on  $s$ .

**Problem 1.** Given the following cost functional and boundary conditions, use the ODEs found in 22.2, 22.4, and 22.5 to solve for and plot the optimal path (position as a function of time,  $x(t)$ ) and acceleration (control as a function of time,  $u(t) = \ddot{x}(t)$ ).

$$\begin{aligned} J(u) &= \int_0^{30} x^2 + \frac{2\pi}{5} u^2 dt \\ s(t) &= \begin{bmatrix} x(t) \\ x'(t) \end{bmatrix}, s(0) = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \quad s(30) = \begin{bmatrix} 16 \\ 10 \end{bmatrix} \end{aligned}$$

Plot your solutions for the optimal  $x(t)$  (position) and  $u(t)$  (acceleration).

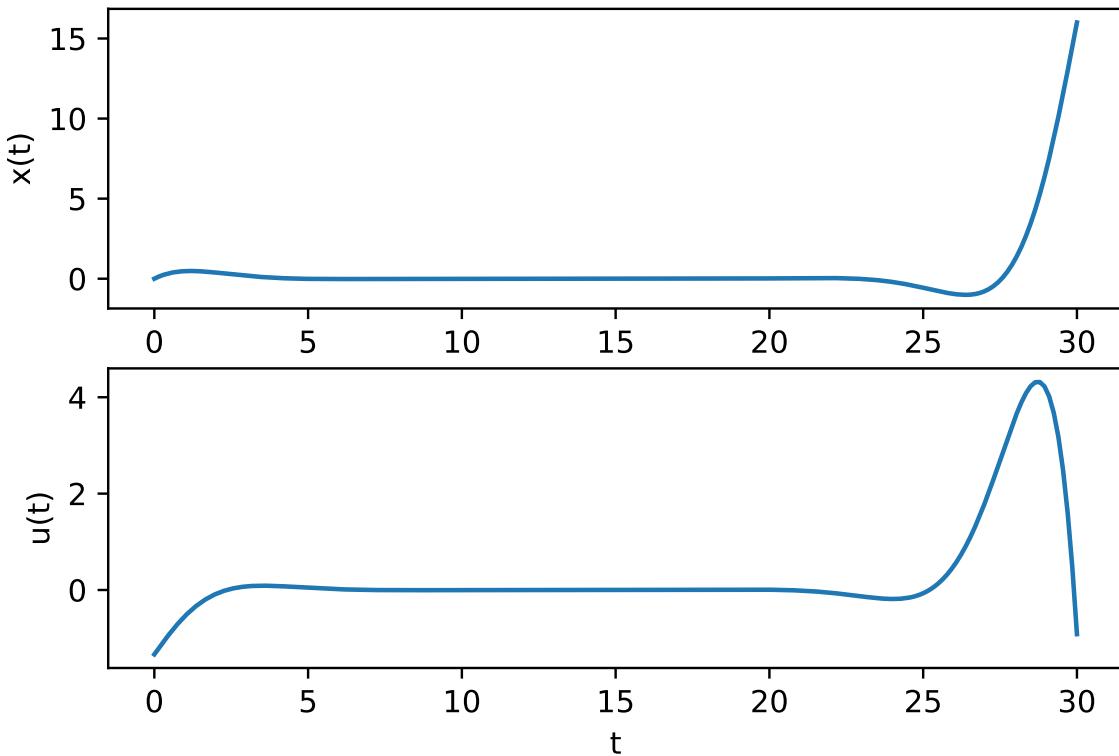


Figure 22.1: Solution to problem 1

## Avoiding Collision

We now expand upon the technique learned above by adding an obstacle in our path. One area of application that relies heavily on optimal control is autonomous driving. A common problem in autonomous driving is the avoidance of obstacles. In this section we will outline a naïve solution to obstacle avoidance with a fixed time horizon.

First we can begin by defining our state variable  $s$ . We will want to understand the position and velocity at a given time so we will define the following state variable

$$s(t) = \begin{bmatrix} x(t) \\ y(t) \\ \dot{x}(t) \\ \dot{y}(t) \end{bmatrix} = \begin{bmatrix} s_1(t) \\ s_2(t) \\ s_3(t) \\ s_4(t) \end{bmatrix}, \quad (22.6)$$

which allows us to track those states in  $\mathbb{R}^2$ .

We can then establish the ODE defined in equation 22.2 by examining  $\dot{s}(t)$

$$\dot{s}(t) = \begin{bmatrix} \dot{s}_1(t) \\ \dot{s}_2(t) \\ \dot{s}_3(t) \\ \dot{s}_4(t) \end{bmatrix} = \begin{bmatrix} \dot{x}(t) \\ \dot{y}(t) \\ \ddot{x}(t) \\ \ddot{y}(t) \end{bmatrix},$$

and if we define our control  $u_1$  and  $u_2$  to be acceleration in the  $x$  and  $y$  directions respectively, then we have

$$\dot{s}(t) = f(t, s, u) = \begin{bmatrix} s_3(t) \\ s_4(t) \\ u_1(t) \\ u_2(t) \end{bmatrix}. \quad (22.7)$$

Next we will define an obstacle. Since we are using integration to define cost, a reasonable way to model an obstacle in this problem would be to use a function. It would be helpful if this function is malleable, allowing us to reposition and resize it, based on the needs of the specific situation. This function also needs to have a large, preferably positive, value in a concentrated location, and it needs to vanish relatively quickly. A decent selection could be a function based on an ellipse, such as this function

$$C(x, y) = \frac{W_1}{((x - c_x)^2/r_x + (y - c_y)^2/r_y)^\lambda + 1}. \quad (22.8)$$

With the function 22.8 we can manipulate the center by changing  $c_x$  and  $c_y$ , and we can control the size by changing  $r_x$  and  $r_y$ . Changing the constant  $W_1$  allows us to change the relative penalty of occupying the same location as the obstacle, and a reasonable value for  $\lambda$  will control the vanishing rate. We will also include a term in the cost functional that weights against high acceleration. This will allow us to model the real world more accurately, though the term we will be using is not a perfect representation of real world acceleration limitations. Our cost functional is the following

$$J(u) = \int_{t_0}^{t_f} 1 + C(x(t), y(t)) + W_2 |u(t)|^2 dt, \quad (22.9)$$

where  $W_2 > 0$  defines the relative penalty of high acceleration. This functional will penalize passing near the obstacle and high levels of acceleration.

With the cost functional defined, we can now create the Hamiltonian and the rest of our BVP. We get the following Hamiltonian

$$H(t, p, s, u) = p_1 s_3 + p_2 s_4 + p_3 u_1 + p_4 u_2 - \left( 1 + C(x, y) + W_2 |u|^2 \right), \quad (22.10)$$

which gives the following costate ODE by equation 22.4

$$\dot{p} = \begin{bmatrix} \dot{p}_1 \\ \dot{p}_2 \\ \dot{p}_3 \\ \dot{p}_4 \end{bmatrix} = \begin{bmatrix} C_x(x, y) \\ C_y(x, y) \\ -p_1 \\ -p_2 \end{bmatrix}. \quad (22.11)$$

Since we're given  $H_u = 0$  in equation 22.5, then we also have the following relations

$$\begin{aligned} u_1(t) &= \frac{1}{2W_2} p_3(t) \\ u_2(t) &= \frac{1}{2W_2} p_4(t). \end{aligned} \quad (22.12)$$

**Problem 2.** Using the ODEs found in 22.7 and 22.11, the obstacle function 22.8, and the following boundary conditions and parameters solve for and plot the optimal path.

$$\begin{aligned} t_0 &= 0, & t_f &= 20 \\ (c_x, c_y) &= (4, 1) \\ (r_x, r_y) &= (5, .5) \\ \lambda &= 20 \\ s_0 &= \begin{bmatrix} 6 \\ 1.5 \\ 0 \\ 0 \end{bmatrix}, & s_f &= \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix} \end{aligned}$$

You will need to choose a  $W_1$  and  $W_2$  which allow the solver to find a valid path. If these parameters are not chosen correctly, the solver may find a path which goes through the obstacle, not around it. Plot the obstacle using `plt.contour()` to be certain path doesn't pass through the obstacle.

Hint: The default for a parameter of `solve_bvp()` called `max_nodes` is not large enough. Try at least `max_nodes = 30000`. You may also find it helpful to use the function `partial` from the module `functools` to preset the parameters for the functions you will be using. If you are struggling to find viable values for  $W_1$  and  $W_2$ , try  $W_1 \in (1, 40)$  and  $W_2 \in (0, 9)$ .

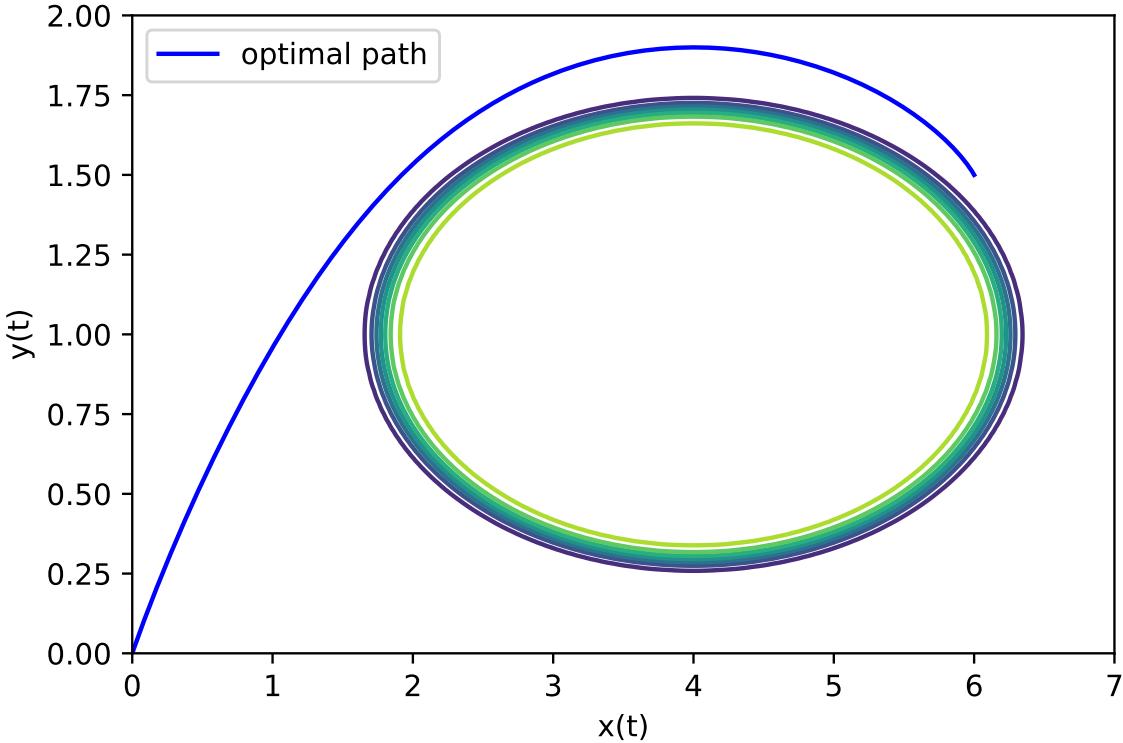


Figure 22.2: Solution to problem 2 for certain choice of parameters. Here we used  $W_1 = 3$  and  $W_2 = 70$ , but those parameters are a choice. Other choices work, too, but will result in a different optimal path around the obstacle.

## Free Time Horizon Problems

In the previous sections and problems, we were working with BVPs that had a fixed start time  $t_0$ , and a fixed end time  $t_f$ . However, we may also encounter systems that have a free end time. In order to solve these problems we will need to make some alterations to the problem. First we will perform a change of basis so that we can work with a fixed end time. Consider the following system

$$\dot{x}(t) = f(x(t), t) \quad t \in [0, t_f],$$

we can do the following change of basis for the time variable

$$\begin{aligned} t &= t_f \hat{t} \\ \implies \frac{d}{dt} &= \frac{d}{d\hat{t}} \frac{dt}{d\hat{t}} \\ \implies \frac{d}{dt} &= \frac{d}{d\hat{t}} \frac{1}{t_f}. \end{aligned}$$

We can now define  $z(\hat{t}) := x(t_f \hat{t})$  which gives us the following new system

$$\dot{z}(\hat{t}) = t_f f(z(\hat{t}), \hat{t}) \quad \hat{t} \in [0, 1].$$

This system can be solved in the same way we solve the fixed time horizon problems. But you may notice that we now have an extra unknown parameter, the final time. Because of this, a free time horizon problem will need one more boundary value to make the system solvable.

So lets examine the earlier example as a free time horizon problem. We start with the ODE system we derived from the second order equation, replacing the fixed final time with a free final time and including the needed third boundary condition

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}' = \begin{bmatrix} y_2 \\ \cos(t) - 9y_1 \end{bmatrix}, \quad y_1(0) = 5/3, \quad y_2(0) = 5, \quad y_1(t_f) = -\frac{5}{3}.$$

Now we make the coordinate change giving the following system

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}' = t_f \begin{bmatrix} z_2 \\ \cos(\hat{t}) - 9z_1 \end{bmatrix}, \quad z_1(0) = 5/3, \quad z_2(0) = 5, \quad z_1(1) = -\frac{5}{3}. \quad (22.13)$$

Now we can solve this system using `solve_bvp` in python. The new argument `p` that we have included in `ode()` and `bc()` is an `ndarray` that contains our parameter  $t_f$ .

```
def ode(t,y,p):
    ''' define the ode system '''
    return p[0]*np.array([y[1], np.cos(t) - 9*y[0]])

def bc(ya,yb,p):
    ''' define the boundary conditions '''
    return np.array([ya[0] - (5/3), ya[1] - 5, yb[0] + 5/3])

# give the time domain
t_steps = 100
t = np.linspace(0,1,t_steps)

# give an initial guess
y0 = np.ones((2,t_steps))
p0 = np.array([6])

# solve the system
sol = solve_bvp(ode, bc, t, y0, p0)
```

The attribute `sol.p[0]` will give the final time the solver found.

When plotting we need to make sure that we remember that  $x(t_f \hat{t}) = z(\hat{t})$ , so we plot in the following way

```
plt.plot(sol.p[0]*t,sol.sol(t)[0])
plt.xlabel('t')
plt.ylabel('y(t)')
plt.show()
```

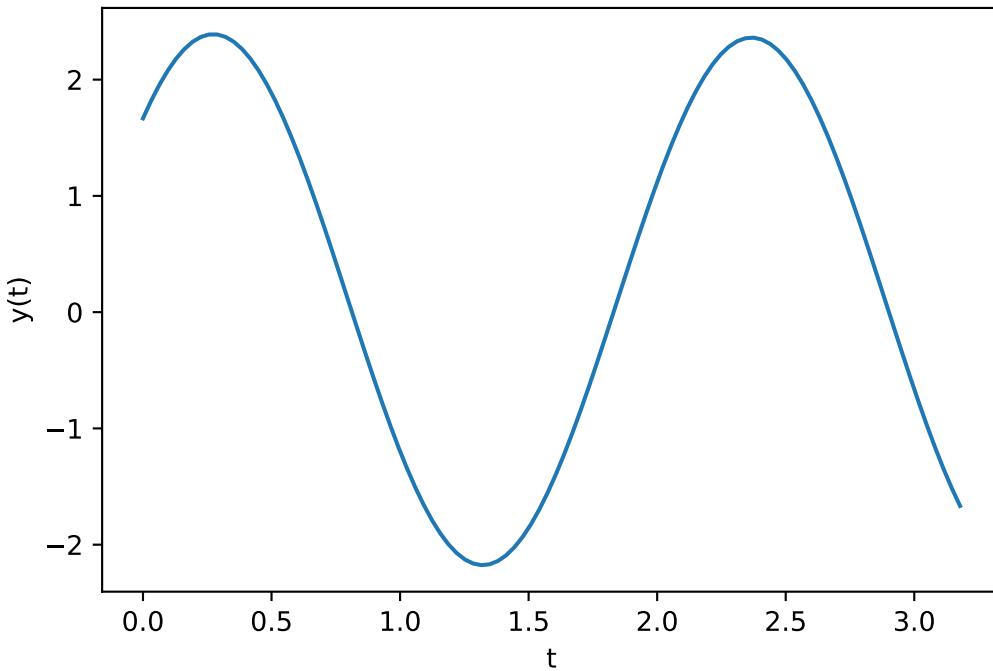


Figure 22.3: The solution to 22.13

**Problem 3.** Solve the following boundary value problem:

$$y'' + 3y = \sin(t)$$

$$y(0) = 0, \quad y(t_f) = \frac{\pi}{2}, \quad y'(t_f) = \frac{1}{2} \left( \sqrt{3}\pi \cot(\pi\sqrt{75}) - 1 \right).$$

Plot your solution. What  $t_f$  did the solver find?

## Free Time, Fixed Endpoint Control Problems

Now that we understand how to formulate free time horizon problems, we can modify our optimal control BVP to become a free time horizon problem. This is actually the best way to formulate many optimal control problems, as we usually don't know exactly how long it takes to traverse the optimal path. The methodology is exactly the same as we used in the last problem, we only need to find the extra boundary value which will allow us to make the end time a free variable.

To find this extra boundary value we will use the fact that the Hamiltonian is 0 for all  $t$  along the optimal path. It is standard to use the final time as the representative so we will assert that

$$H(t_f, p(t_f), s(t_f), u(t_f)) = 0. \quad (22.14)$$

You may notice that when you solve an optimal control problem as a free end time BVP, the optimal path you get is different than what you found when it was a fixed end time BVP. This is because the free end time solution actually arrives faster. The solution found in the fixed end time formulation is the optimal path for a certain fixed end time, but it may not be the overall fastest path that avoids the obstacle.

**Problem 4.** Refactor your code from problem 2 to create a free end time BVP and use a new boundary value derived from 22.14. Plot the solution you found. What is the optimal time?

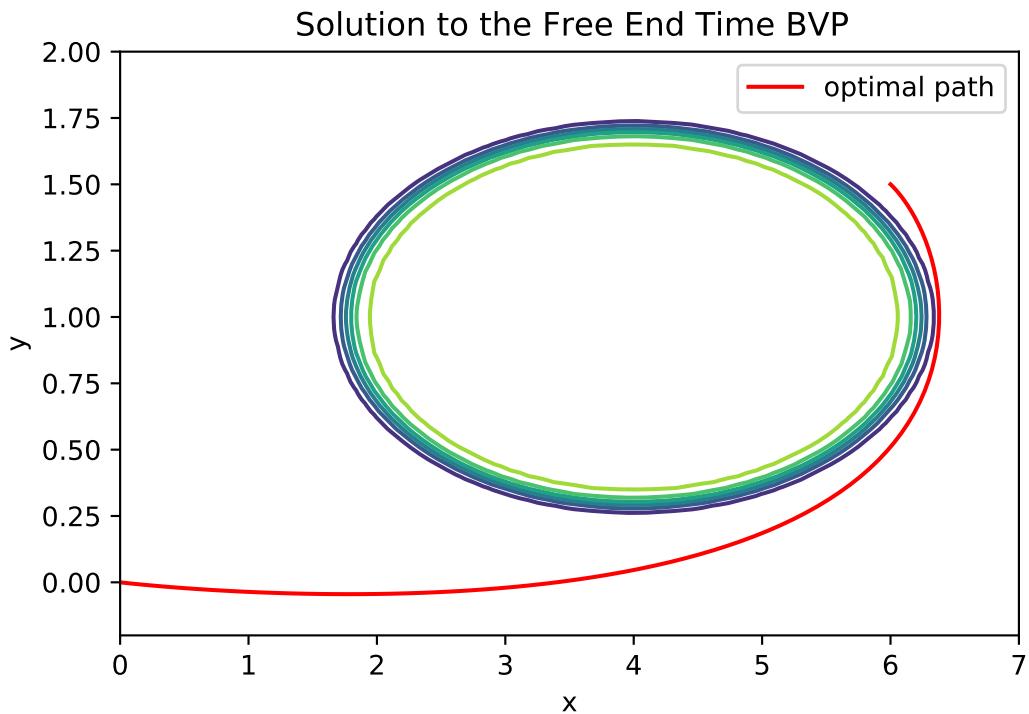


Figure 22.4: The solution to 4. Here we used  $W_1 = 4$  and  $W_2 = 0.1$  and got an optimal time of about 4.47. Those parameters are a choice. Other choices work, too, but will result in a different optimal path around the obstacle and in a different optimal time.

# 23 The Inverted Pendulum

**Lab Objective:** We will set up the LQR optimal control problem for the inverted pendulum and compute the solution numerically.

Think back to your childhood days when, for entertainment purposes, you'd balance objects: a book on your head, a spoon on your nose, or even a broom on your hand. Learning how to walk was likely your initial introduction to the inverted pendulum problem.

A pendulum has two rest points: a stable rest point directly underneath the pivot point of the pendulum, and an unstable rest point directly above. The generic pendulum problem is to simply describe the dynamics of the object on the pendulum (called the ‘bob’). The inverted pendulum problem seeks to guide the bob toward the unstable fixed point at the top of the pendulum. Since the fixed point is unstable, the bob must be balanced relentlessly to keep it upright.

The inverted pendulum is an important classical problem in dynamics and control theory, and is often used to test different control strategies. One application of the inverted pendulum is the guidance of rockets and missiles. Aerodynamic instability occurs because the center of mass of the rocket is not the same as the center of drag. Small gusts of wind or variations in thrust require constant attention to the orientation of the rocket.

## The Simple Pendulum

We begin by studying the simple pendulum setting. Suppose we have a pendulum consisting of a bob with mass  $m$  rotating about a pivot point at the end of a (massless) rod of length  $l$ . Let  $\theta(t)$  represent the angular displacement of the bob from its stable equilibrium. By Hamilton’s Principle, the path  $\theta$  that is taken by the bob minimizes the functional

$$J[\theta] = \int_{t_0}^{t_1} L, \quad (23.1)$$

where the Lagrangian  $L = T - U$  is the difference between the kinetic and potential energies of the bob.

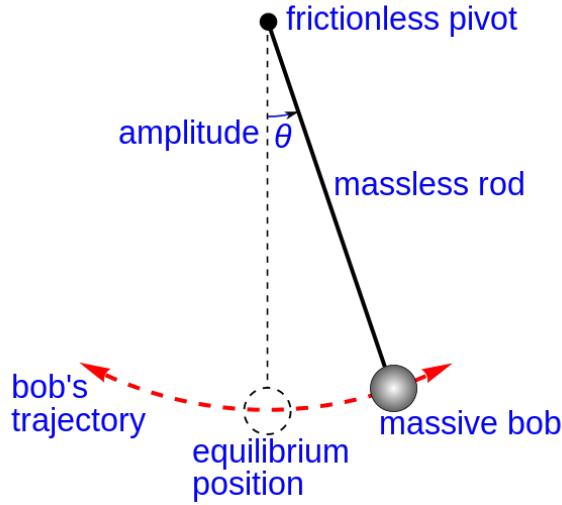


Figure 23.1: The frame of reference for the simple pendulum problem.

The kinetic energy of the bob is given by  $mv^2/2$ , where  $v$  is the velocity of the bob. In terms of  $\theta$ , the kinetic energy becomes

$$\begin{aligned} T &= \frac{m}{2}v^2 = \frac{m}{2}(\dot{x}^2 + \dot{y}^2), \\ &= \frac{m}{2}((l \cos(\theta)\dot{\theta})^2 + (l \sin(\theta)\dot{\theta})^2), \\ &= \frac{ml^2\dot{\theta}^2}{2}. \end{aligned} \tag{23.2}$$

The potential energy of the bob is  $U = mg(l - l \cos \theta)$ . From these expressions we can form the Euler-Lagrange equation, which determines the path of the bob:

$$\begin{aligned} 0 &= L_\theta - \frac{d}{dx}L_{\dot{\theta}}, \\ &= -mgl \sin \theta - ml^2\ddot{\theta}, \\ &= \ddot{\theta} + \frac{g}{l} \sin \theta. \end{aligned} \tag{23.3}$$

Since in this setting the energy of the pendulum is conserved, the equilibrium position  $\theta = 0$  is only Lyapunov stable. When forces such as friction and air drag are considered  $\theta = 0$  becomes an asymptotically stable equilibrium.

## The Inverted Pendulum

### The Control System

We consider a gift suspended above a rickshaw by a (massless) rod of length  $l$ . The rickshaw and its suspended gift will have masses  $M$  and  $m$  respectively,  $M > m$ . Let  $\theta$  represent the angle between the gift and its unstable equilibrium, with clockwise orientation. Let  $v_1$  and  $v_2$  represent the velocities of the rickshaw and the gift, and  $F$  the force exerted on the rickshaw. The rickshaw will be restricted to traveling along a straight line (the  $x$ -axis).

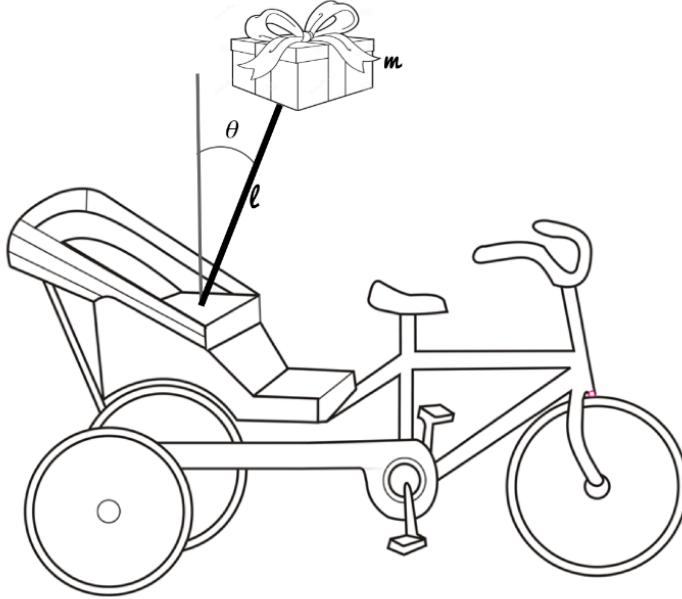


Figure 23.2: The inverted pendulum problem on a mobile rickshaw with a present suspended above.

By Hamilton's Principle, the path  $(x, \theta)$  of the rickshaw and the present minimizes the functional

$$J[x, \theta] = \int_{t_0}^{t_1} L, \quad (23.4)$$

where the Lagrangian  $L = T - U$  is the difference between the kinetic energy of the present on the pendulum, and its potential energy.

Since the position of the rickshaw and the present are  $(x(t), 0)$  and  $(x - l \sin \theta, l \cos \theta)$  respectively, the total kinetic energy is

$$\begin{aligned} T &= \frac{1}{2} M v_1^2 + \frac{1}{2} m v_2^2, \\ &= \frac{1}{2} M \dot{x}^2 + \frac{1}{2} m ((\dot{x} - l \dot{\theta} \cos \theta)^2 + (-l \dot{\theta} \sin \theta)^2), \\ &= \frac{1}{2} (M + m) \dot{x}^2 + \frac{1}{2} m l^2 \dot{\theta}^2 - m l \dot{x} \dot{\theta} \cos \theta. \end{aligned} \quad (23.5)$$

The total potential energy is

$$U = m g l \cos \theta.$$

The path  $(x, \theta)$  of the rickshaw and the present satisfy the Euler-Lagrange differential equations, but the problem involves a nonconservative force  $F$  acting in the  $x$  direction. By way of D'Alambert's Principle, our normal Euler-Lagrange equations now include the nonconservative force  $F$  on the right side of the equation:

$$\begin{aligned} \frac{\partial L}{\partial x} - \frac{d}{dt} \frac{\partial L}{\partial \dot{x}} &= F, \\ \frac{\partial L}{\partial \theta} - \frac{d}{dt} \frac{\partial L}{\partial \dot{\theta}} &= 0. \end{aligned} \quad (23.6)$$

After expanding (23.6) we see that  $x(t)$  and  $\theta(t)$  satisfy

$$\begin{aligned} F &= (M+m)\ddot{x} - ml\ddot{\theta} \cos \theta + ml\dot{\theta}^2 \sin \theta, \\ l\ddot{\theta} &= g \sin \theta + \ddot{x} \cos \theta. \end{aligned} \quad (23.7)$$

At this point we make a further simplifying assumption. If  $\theta$  starts close to 0, we may assume that the corresponding force  $F$  will keep  $\theta$  small. In this case, we linearize (23.7) about  $(\theta, \dot{\theta}) = (0, 0)$ , obtaining the equations

$$\begin{aligned} F &= (M+m)\ddot{x} - ml\ddot{\theta}, \\ l\ddot{\theta} &= g\theta + \ddot{x}. \end{aligned}$$

These equations can be further manipulated to obtain

$$\begin{aligned} \ddot{x} &= \frac{1}{M}F - \frac{m}{M}g\theta, \\ \ddot{\theta} &= \frac{1}{Ml}F + \frac{g}{Ml}(M+m)\theta. \end{aligned} \quad (23.8)$$

We will now write (23.8) as a first order system. Making the assignments  $x_1 = x$ ,  $x_2 = x'_1$ ,  $\theta_1 = \theta$ ,  $\theta_2 = \theta'_1$ , letting  $u = F$  represent the control variable, we obtain

$$\begin{bmatrix} x_1 \\ x_2 \\ \theta_1 \\ \theta_2 \end{bmatrix}' = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & \frac{mg}{M} & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & \frac{g}{Ml}(M+m) & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \theta_1 \\ \theta_2 \end{bmatrix} + u \begin{bmatrix} 0 \\ \frac{1}{M} \\ 0 \\ \frac{1}{Ml} \end{bmatrix},$$

which can be written more concisely as

$$z' = Az + Bu.$$

### The infinite time horizon LQR problem

We consider the cost function

$$\begin{aligned} J[z] &= \int_0^\infty (q_1x_1^2 + q_2x_2^2 + q_3\theta_1^2 + q_4\theta_2^2 + ru^2) dt \\ &= \int_0^\infty z^T Q z + u^T R u dt \end{aligned} \quad (23.9)$$

where  $q_1, q_2, q_3, q_4$ , and  $r$  are nonnegative weights, and

$$Q = \begin{bmatrix} q_1 & 0 & 0 & 0 \\ 0 & q_2 & 0 & 0 \\ 0 & 0 & q_3 & 0 \\ 0 & 0 & 0 & q_4 \end{bmatrix}, R = [r].$$

**Problem 1.** Write a function that returns the matrices  $A, B, Q$ , and  $R$  given above. Let  $g = 9.8 \text{ m/s}^2$ .

```
def linearized_init(M, m, l, q1, q2, q3, q4, r):
    ...
Parameters:
```

```

-----
M, m: floats
    masses of the rickshaw and the present
l   : float
    length of the rod
q1, q2, q3, q4, r : floats
    relative weights of the position and velocity of the rickshaw, ←
        the
    angular displacement theta and the change in theta, and the ←
        control

Return
-----
A : ndarray of shape (4,4)
B : ndarray of shape (4,1)
Q : ndarray of shape (4,4)
R : ndarray of shape (1,1)
...
pass

```

The optimal control problem (23.9) is an example of a Linear Quadratic Regulator (LQR), and is known to have an optimal control  $\tilde{u}$  described by a linear state feedback law:

$$\tilde{u} = -R^{-1}B^T P \tilde{z}.$$

Here  $P$  is a matrix function that satisfies the Riccati differential equation (RDE)

$$\dot{P}(t) = PA + A^T P + Q - PBR^{-1}B^T P.$$

Since this problem has an infinite time horizon, we have  $\dot{P} = 0$ . Thus  $P$  is a constant matrix, and can be found by solving the algebraic Riccati equation (ARE)

$$PA + A^T P + Q - PBR^{-1}B^T P = 0. \quad (23.10)$$

The evolution of the optimal state vector  $\tilde{z}$  can then be described by <sup>1</sup>

$$\dot{\tilde{z}} = (A - BR^{-1}B^T P)\tilde{z}. \quad (23.11)$$

**Problem 2.** Write the following function to find the matrix  $P$ . Use `scipy.optimize.root`. Since `root` takes in a vector and not a matrix, you will have to reshape the matrix  $P$  before passing it in and after getting your result, using `np.reshape(16)` and `np.reshape((4,4))`.

```

def find_P(A, B, Q, R):
    ...
Parameters:
-----
A, Q      : ndarrays of shape (4,4)

```

<sup>1</sup>See Calculus of Variations and Optimal Control Theory, Daniel Liberzon, Ch.6

```

B      : ndarray of shape (4,1)
R      : ndarray of shape (1,1)

Returns
-----
P      : the matrix solution of the Riccati equation
...
pass

```

Using the values

```

M, m = 23., 5.
l = 4.
q1, q2, q3, q4 = 1., 1., 1., 1.
r = 10.

```

compute the eigenvalues of  $A - BR^{-1}B^T P$ . Are any of the eigenvalues positive? Consider differential equation (23.11) governing the optimal state  $\tilde{z}$ . Using this value of  $P$ , will we necessarily have  $\tilde{z} \rightarrow 0$ ?

**Problem 3.** Write the following function that implements the LQR solution described earlier. Use `scipy.integrate.solve_ivp` to solve the IVP.

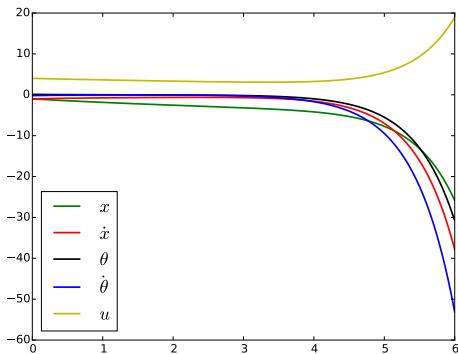
```

def rickshaw(tv, X0, A, B, Q, R, P):
    ...
    Parameters:
    -----
    tv   : ndarray of time values, with shape (n+1,)
    X0  : Initial conditions on state variables
    A, Q: ndarrays of shape (4,4)
    B   : ndarray of shape (4,1)
    R   : ndarray of shape (1,1)
    P   : ndarray of shape (4,4)

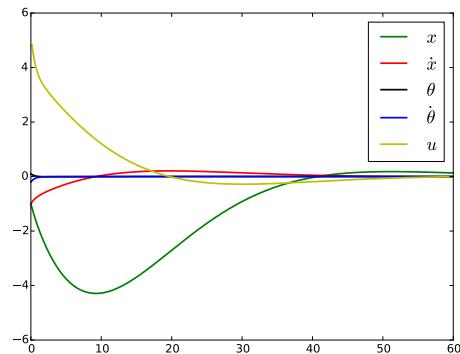
    Returns
    -----
    Z : ndarray of shape (n+1,4), the state vector at each time
    U : ndarray of shape (n+1,), the control values
    ...
    pass

```

Notice that we have no information on how many solutions (23.10) possesses. In general there may be many solutions. We hope to find a unique solution  $P$  that is stabilizing: the eigenvalues of  $A - BR^{-1}B^T P$  have negative real part. To find this  $P$ , use the function `solve_continuous_are` from `scipy.linalg`. This function is designed to solve the continuous algebraic Riccati equation.



$P$  is found using `scipy.optimize.root`.



$P$  is found using `solve_continuous_are`.

Figure 23.3: The solutions of Problem 4.

**Problem 4.** Test the function made in Problem (3) with the following inputs:

```
M, m = 23., 5.
l = 4.
q1, q2, q3, q4 = 1., 1., 1., 1.
r = 10.
tf = None
X0 = np.array([-1, -1, .1, -.2])
```

Find the matrix  $P$  using the `scipy.optimize.root` method with  $tf=6$  as well as the `solve_continuous_are` method with  $tf=60$ . Plot the solutions  $\tilde{z}$  and  $\tilde{u}$ . Compare your results as shown in Figure 23.3.



## Part II

# Appendices



# A

## Getting Started

The labs in this curriculum aim to introduce computational and mathematical concepts, walk through implementations of those concepts in Python, and use industrial-grade code to solve interesting, relevant problems. Lab assignments are usually about 5–10 pages long and include code examples (yellow boxes), important notes (green boxes), warnings about common errors (red boxes), and about 3–7 exercises (blue boxes). Get started by downloading the lab manual(s) for your course from <http://foundations-of-applied-mathematics.github.io/>.

### Submitting Assignments

#### Labs

Every lab has a corresponding specifications file with some code to get you started and to make your submission compatible with automated test drivers. Like the lab manuals, these materials are hosted at <http://foundations-of-applied-mathematics.github.io/>.

Download the `.zip` file for your course, unzip the folder, and move it somewhere where it won't get lost. This folder has some setup scripts and a collection of folders, one per lab, each of which contains the specifications file(s) for that lab. See [Student-Materials/wiki/Lab-Index](#) for the complete list of labs, their specifications and data files, and the manual that each lab belongs to.

#### Achtung!

Do **not** move or rename the lab folders or the enclosed specifications files; if you do, the test drivers will not be able to find your assignment. Make sure your folder and file names match [Student-Materials/wiki/Lab-Index](#).

To submit a lab, modify the provided specifications file and use the file-sharing program specified by your instructor (discussed in the next section). The instructor will drop feedback files in the lab folder after grading the assignment. For example, the Introduction to Python lab has the specifications file `PythonIntro/python_intro.py`. To complete that assignment, modify `PythonIntro/python_intro.py` and submit it via your instructor's file-sharing system. After grading, the instructor will create a file called `PythonIntro/PythonIntro_feedback.txt` with your score and some feedback.

## Homework

Non-lab coding homework should be placed in the `_Homework/` folder and submitted like a lab assignment. Be careful to name your assignment correctly so the instructor (and test driver) can find it. The instructor may drop specifications files and/or feedback files in this folder as well.

## Setup

### Achtung!

We strongly recommend using a Unix-based operating system (Mac or Linux) for the labs. Unix has a true bash terminal, works well with git and python, and is the preferred platform for computational and data scientists. It is possible to do this curriculum with Windows, but expect some road bumps along the way.

There are two ways to submit code to the instructor: with git (<http://git-scm.com/>), or with a file-syncing service like Google Drive. Your instructor will indicate which system to use.

### Setup With Git

Git is a program that manages updates between an online code repository and the copies of the repository, called clones, stored locally on computers. If git is not already installed on your computer, download it at <http://git-scm.com/downloads>. If you have never used git, you might want to read a few of the following resources.

- Official git tutorial: <https://git-scm.com/docs/gittutorial>
- Bitbucket git tutorials: <https://www.atlassian.com/git/tutorials>
- GitHub git cheat sheet: [services.github.com/.../github-git-cheat-sheet.pdf](https://services.github.com/.../github-git-cheat-sheet.pdf)
- GitLab git tutorial: <https://docs.gitlab.com/ce/gitlab-basics/start-using-git.html>
- Codecademy git lesson: <https://www.codecademy.com/learn/learn-git>
- Training video series by GitHub: <https://www.youtube.com/playlist?list=PLg7.../>

There are many websites for hosting online git repositories. Your instructor will indicate which web service to use, but we only include instructions here for setup with Bitbucket.

1. Sign up. Create a Bitbucket account at <https://bitbucket.org>. If you use an academic email address (ending in `.edu`, etc.), you will get free unlimited public and private repositories.
2. Make a new repository. On the Bitbucket page, click the `+` button from the menu on the left and, under **CREATE**, select **Repository**. Provide a name for the repository, mark the repository as **private**, and make sure the repository type is **Git**. For **Include a README?**, select **No** (if you accidentally include a README, delete the repository and start over). Under **Advanced settings**, enter a short description for your repository, select **No forks** under forking, and select **Python** as the language. Finally, click the blue **Create repository** button. Take note of the URL of the webpage that is created; it should be something like <https://bitbucket.org/<name>/<repo>>.

3. Give the instructor access to your repository. On your newly created Bitbucket repository page (<https://bitbucket.org/<name>/<repo>> or similar), go to **Settings** in the menu to the left and select **User and group access**, the second option from the top. Enter your instructor's Bitbucket username under **Users** and click **Add**. Select the blue **Write** button so your instructor can read from and write feedback to your repository.
4. Create an SSH key. This step needs to be done only once on each computer that you want to be able to use to access your repository. If you have multiple repositories on the same computer, you do not need to repeat this step for each one. To create an SSH key, in a shell application (Terminal on Linux or Mac, or Git Bash (<https://gitforwindows.org/>) on Windows), enter the following command:

```
$ ssh-keygen
```

Press the Enter or Return key to accept the default file location. It will then prompt to enter a passphrase; this acts as a password to use the SSH key. If you do not want a passphrase, leave it blank and press Enter again. The key will then be created. The file for the key will be placed in in the `/home/<username>/ .ssh` directory on Linux; in `/Users/<username>/ .ssh` on macOS; and in `/c/users/<username>/ .ssh` on Windows.

Now that the key is created, you need to add it to your Bitbucket account. From Bitbucket, choose **Personal settings** and then **SSH keys**. Click **Add key** and enter a label (what it is doesn't matter). Now, using the file explorer, navigate to the SSH key you created, and open the public key file. The file will be called something like `id_rsa.pub`; do NOT use `id_rsa` (without the `.pub` extension). Copy the contents of this file, paste it into the Key field on Bitbucket, and press Save.

For more options and some troubleshooting information, refer to <https://support.atlassian.com/bitbucket-cloud/docs/set-up-an-ssh-key/>.

5. Connect your folder to the new repository. In a shell application (Terminal on Linux or Mac, or Git Bash (<https://gitforwindows.org/>) on Windows), enter the following commands.

```
# Navigate to your folder.
$ cd /path/to/folder # cd means 'change directory'.


# Make sure you are in the right place.
$ pwd # pwd means 'print working directory'.
/path/to/folder
$ ls *.md # ls means 'list files'.
README.md # This means README.md is in the working directory.


# Connect this folder to the online repository.
$ git init
$ git remote add origin git@bitbucket.org:<name>/<repo>.git


# Record your credentials.
$ git config --local user.name "your name"
$ git config --local user.email "your email"


# Add the contents of this folder to git and update the repository.
```

```
$ git add --all
$ git commit -m "initial commit"
$ git push origin master
```

For example, if your Bitbucket username is `greek314`, the repository is called `acmev1`, and the folder is called `Student-Materials/` and is on the desktop, enter the following commands.

```
# Navigate to the folder.
$ cd ~/Desktop/Student-Materials

# Make sure this is the right place.
$ pwd
/Users/Archimedes/Desktop/Student-Materials
$ ls *.md
README.md

# Connect this folder to the online repository.
$ git init
$ git remote add origin git@bitbucket.org:greek314/acmev1.git

# Record credentials.
$ git config --local user.name "archimedes"
$ git config --local user.email "greek314@example.com"

# Add the contents of this folder to git and update the repository.
$ git add --all
$ git commit -m "initial commit"
$ git push origin master
```

At this point you should be able to see the files on your repository page from a web browser. If you enter the repository URL incorrectly in the `git remote add origin` step, you can reset it with the following line:

```
$ git remote set-url origin git@bitbucket.org:<name>/<repo>.git
```

### Note

You may get the an error like the following when you run `git push`:

```
remote: Bitbucket Cloud recently stopped supporting account passwords←
      for Git authentication.
...
fatal: Authentication failed for 'https://bitbucket.org/<name>/<repo←
>.git/'
```

If this error occurs, your repository URL is in the wrong format; most likely, you used the `https` version instead of what is shown above. You can use the `git remote set-url origin` command to fix this issue as well.

6. Download data files. Many labs have accompanying data files. To download these files, navigate to your clone and run the `download_data.sh` bash script, which downloads the files and places them in the correct lab folder for you. You can also find individual data files through [Student-Materials/wiki/Lab-Index](#).

```
# Navigate to your folder and run the script.
$ cd /path/to/folder
$ bash download_data.sh
```

7. Install Python package dependencies. The labs require several third-party Python packages that don't come bundled with Anaconda. Run the following command to install the necessary packages.

```
# Navigate to your folder and run the script.
$ cd /path/to/folder
$ bash install_dependencies.sh
```

8. (Optional) Clone your repository. If you want your repository on another computer after completing steps 1–5, use the following commands.

```
# Navigate to where you want to put the folder.
$ cd ~/Desktop/or/something/

# Clone the folder from the online repository.
$ git clone git@bitbucket.org:<name>/<repo>.git <foldername>

# Record your credentials in the new folder.
$ cd <foldername>
$ git config --local user.name "your name"
$ git config --local user.email "your email"

# Download data files to the new folder.
$ bash download_data.sh
```

## Setup Without Git

Even if you aren't using git to submit files, you must install it (<http://git-scm.com/downloads>) in order to get the data files for each lab. Share your folder with your instructor according to their directions, and follow steps 6 and 7 of the previous section to download the data files and install package dependencies.

## Using Git

Git manages the history of a file system through commits, or checkpoints. Use `git status` to see the files that have been changed since the last commit. These changes are then moved to the staging area, a list of files to save during the next commit, with `git add <filename(s)>`. Save the changes in the staging area with `git commit -m "<A brief message describing the changes>"`.

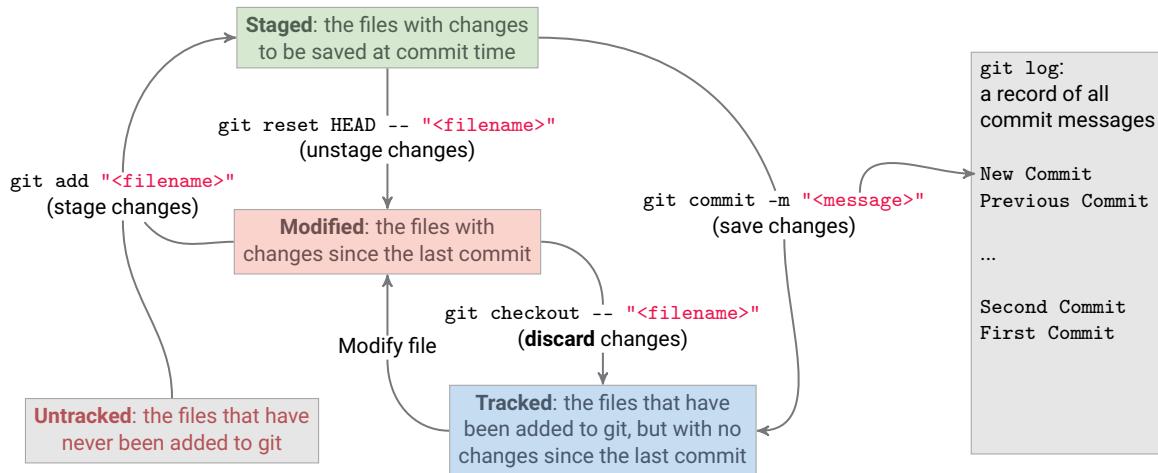


Figure A.1: Git commands to stage, unstage, save, or discard changes. Commit messages are recorded in the log.

All of these commands are done within a clone of the repository, stored somewhere on a computer. This repository must be manually synchronized with the online repository via two other git commands: `git pull origin master`, to pull updates from the web to the computer; and `git push origin master`, to push updates from the computer to the web.

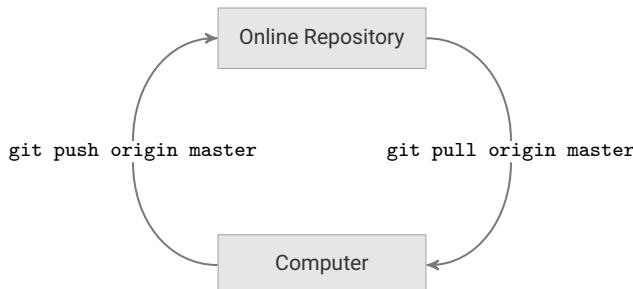


Figure A.2: Exchanging git commits between the repository and a local clone.

Command	Explanation
<code>git status</code>	Display the staging area and untracked changes.
<code>git pull origin master</code>	Pull changes from the online repository.
<code>git push origin master</code>	Push changes to the online repository.
<code>git add &lt;filename(s)&gt;</code>	Add a file or files to the staging area.
<code>git add -u</code>	Add all modified, tracked files to the staging area.
<code>git commit -m "&lt;message&gt;"</code>	Save the changes in the staging area with a given message.
<code>git checkout -- &lt;filename&gt;</code>	Revert changes to an unstaged file since the last commit.
<code>git reset HEAD -- &lt;filename&gt;</code>	Remove a file from the staging area.
<code>git diff &lt;filename&gt;</code>	See the changes to an unstaged file since the last commit.
<code>git diff --cached &lt;filename&gt;</code>	See the changes to a staged file since the last commit.
<code>git config --local &lt;option&gt;</code>	Record your credentials ( <code>user.name</code> , <code>user.email</code> , etc.).

Table A.1: Common git commands.

### Note

When pulling updates with `git pull origin master`, your terminal may sometimes display the following message.

```
Merge branch 'master' of git@bitbucket.org:<name>/<repo> into master

# Please enter a commit message to explain why this merge is necessary,
# especially if it merges an updated upstream into a topic branch.
#
# Lines starting with '#' will be ignored, and an empty message aborts
# the commit.
~
```

This means that someone else (the instructor) has pushed a commit that you do not yet have, while you have also made one or more commits locally that they do not have. This screen, displayed in vim ([https://en.wikipedia.org/wiki/Vim\\_\(text\\_editor\)](https://en.wikipedia.org/wiki/Vim_(text_editor))), is asking you to enter a message (or use the default message) to create a merge commit that will reconcile both changes. To close this screen and create the merge commit, type :wq and press `enter`.

## Example Work Sessions

```
$ cd ~/Desktop/Student-Materials/
$ git pull origin master                                # Pull updates.
### Make changes to a file.
$ git add -u                                            # Track changes.
$ git commit -m "Made some changes."                   # Commit changes.
$ git push origin master                               # Push updates.
```

```
# Pull any updates from the online repository (such as TA feedback).
$ cd ~/Desktop/Student-Materials/
$ git pull origin master
From bitbucket.org:username/repo
 * branch            master      -> FETCH_HEAD
Already up-to-date.

### Work on the labs. For example, modify PythonIntro/python_intro.py.

$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
Changes not staged for commit:
  (use "git add <file>..." to update what will be committed)
  (use "git checkout -- <file>..." to discard changes in working directory)

    PythonIntro/python_intro.py

# Track the changes with git.
$ git add PythonIntro/python_intro.py
$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
Changes to be committed:
  (use "git reset HEAD <file>..." to unstage)

    modified:   PythonIntro/python_intro.py

# Commit the changes to the repository with an informative message.
$ git commit -m "Made some changes"
[master fed9b34] Made some changes
  1 file changed, 10 insertion(+) 1 deletion(-)

# Push the changes to the online repository.
$ git push origin master
Counting objects: 3, done.
Delta compression using up to 2 threads.
Compressing objects: 100% (2/2), done.
Writing objects: 100% (3/3), 327 bytes | 0 bytes/s, done.
Total 3 (delta 0), reused 0 (delta 0)
To git@bitbucket.org:username/repo.git
  5742a1b..fed9b34  master -> master

$ git status
On branch master
Your branch is up-to-date with 'origin/master'.
nothing to commit, working directory clean
```

# B

# Installing and Managing Python

**Lab Objective:** One of the great advantages of Python is its lack of overhead: it is relatively easy to download, install, start up, and execute. This appendix introduces tools for installing and updating specific packages and gives an overview of possible environments for working efficiently in Python.

## Installing Python via Anaconda

A Python distribution is a single download containing everything needed to install and run Python, together with some common packages. For this curriculum, we **strongly** recommend using the Anaconda distribution to install Python. Anaconda includes IPython, a few other tools for developing in Python, and a large selection of packages that are common in applied mathematics, numerical computing, and data science. Anaconda is free and available for Windows, Mac, and Linux.

Follow these steps to install Anaconda.

1. Go to <https://www.anaconda.com/download/>.
2. Download the **Python 3.6** graphical installer specific to your machine.
3. Open the downloaded file and proceed with the default configurations.

For help with installation, see <https://docs.anaconda.com/anaconda/install/>. This page contains links to detailed step-by-step installation instructions for each operating system, as well as information for updating and uninstalling Anaconda.

### Achtung!

This curriculum uses Python 3.6, **not** Python 2.7. With the wrong version of Python, some example code within the labs may not execute as intended or result in an error.

## Managing Packages

A Python package manager is a tool for installing or updating Python packages, which involves downloading the right source code files, placing those files in the correct location on the machine, and linking the files to the Python interpreter. **Never** try to install a Python package without using a package manager (see <https://xkcd.com/349/>).

### Conda

Many packages are not included in the default Anaconda download but can be installed via Anaconda's package manager, `conda`. See <https://docs.anaconda.com/anaconda/packages/pkg-docs> for the complete list of available packages. When you need to update or install a package, **always** try using `conda` first.

Command	Description
<code>conda install &lt;package-name&gt;</code>	Install the specified package.
<code>conda update &lt;package-name&gt;</code>	Update the specified package.
<code>conda update conda</code>	Update <code>conda</code> itself.
<code>conda update anaconda</code>	Update <b>all</b> packages included in Anaconda.
<code>conda --help</code>	Display the documentation for <code>conda</code> .

For example, the following terminal commands attempt to install and update `matplotlib`.

```
$ conda update conda          # Make sure that conda is up to date.
$ conda install matplotlib    # Attempt to install matplotlib.
$ conda update matplotlib     # Attempt to update matplotlib.
```

See <https://conda.io/docs/user-guide/tasks/manage-pkgs.html> for more examples.

#### Note

The best way to ensure a package has been installed correctly is to try importing it in IPython.

```
# Start IPython from the command line.
$ ipython
IPython 6.5.0 -- An enhanced Interactive Python. Type '?' for help.

# Try to import matplotlib.
In [1]: from matplotlib import pyplot as plt      # Success!
```

#### Achtung!

Be careful not to attempt to update a Python package while it is in use. It is safest to update packages while the Python interpreter is not running.

## Pip

The most generic Python package manager is called `pip`. While it has a larger package list, `conda` is the cleaner and safer option. Only use `pip` to manage packages that are not available through `conda`.

Command	Description
<code>pip install package-name</code>	Install the specified package.
<code>pip install --upgrade package-name</code>	Update the specified package.
<code>pip freeze</code>	Display the version number on all installed packages.
<code>pip --help</code>	Display the documentation for <code>pip</code> .

See [https://pip.pypa.io/en/stable/user\\_guide/](https://pip.pypa.io/en/stable/user_guide/) for more complete documentation.

## Workflows

There are several different ways to write and execute programs in Python. Try a variety of workflows to find what works best for you.

### Text Editor + Terminal

The most basic way of developing in Python is to write code in a text editor, then run it using either the Python or IPython interpreter in the terminal.

There are many different text editors available for code development. Many text editors are designed specifically for computer programming which contain features such as syntax highlighting and error detection, and are highly customizable. Try installing and using some of the popular text editors listed below.

- Atom: <https://atom.io/>
- Sublime Text: <https://www.sublimetext.com/>
- Notepad++ (Windows): <https://notepad-plus-plus.org/>
- Geany: <https://www.geany.org/>
- Vim: <https://www.vim.org/>
- Emacs: <https://www.gnu.org/software/emacs/>

Once Python code has been written in a text editor and saved to a file, that file can be executed in the terminal or command line.

```
$ ls                               # List the files in the current directory.
hello_world.py
$ cat hello_world.py               # Print the contents of the file to the terminal.
print("hello, world!")
$ python hello_world.py            # Execute the file.
hello, world!

# Alternatively, start IPython and run the file.
$ ipython
```

```
IPython 6.5.0 -- An enhanced Interactive Python. Type '?' for help.  
  
In [1]: %run hello_world.py  
hello, world!
```

IPython is an enhanced version of Python that is more user-friendly and interactive. It has many features that cater to productivity such as tab completion and object introspection.

### Note

While Mac and Linux computers come with a built-in bash terminal, Windows computers do not. Windows does come with Powershell, a terminal-like application, but some commands in Powershell are different than their bash analogs, and some bash commands are missing from Powershell altogether. There are two good alternatives to the bash terminal for Windows:

- Windows subsystem for linux: [docs.microsoft.com/en-us/windows/wsl/](https://docs.microsoft.com/en-us/windows/wsl/).
- Git bash: <https://gitforwindows.org/>.

## Jupyter Notebook

The Jupyter Notebook (previously known as IPython Notebook) is a browser-based interface for Python that comes included as part of the Anaconda Python Distribution. It has an interface similar to the IPython interpreter, except that input is stored in cells and can be modified and re-evaluated as desired. See <https://github.com/jupyter/jupyter/wiki/> for some examples.

To begin using Jupyter Notebook, run the command `jupyter notebook` in the terminal. This will open your file system in a web browser in the Jupyter framework. To create a Jupyter Notebook, click the **New** drop down menu and choose **Python 3** under the **Notebooks** heading. A new tab will open with a new Jupyter Notebook.

Jupyter Notebooks differ from other forms of Python development in that notebook files contain not only the raw Python code, but also formatting information. As such, Jupyter Notebook files cannot be run in any other development environment. They also have the file extension `.ipynb` rather than the standard Python extension `.py`.

Jupyter Notebooks also support Markdown—a simple text formatting language—and L<sup>A</sup>T<sub>E</sub>X, and can embed images, sound clips, videos, and more. This makes Jupyter Notebook the ideal platform for presenting code.

## Integrated Development Environments

An integrated development environment (IDEs) is a program that provides a comprehensive environment with the tools necessary for development, all combined into a single application. Most IDEs have many tightly integrated tools that are easily accessible, but come with more overhead than a plain text editor. Consider trying out each of the following IDEs.

- JupyterLab: <http://jupyterlab.readthedocs.io/en/stable/>
- PyCharm: <https://www.jetbrains.com/pycharm/>

- Spyder: <http://code.google.com/p/spyderlib/>
- Eclipse with PyDev: <http://www.eclipse.org/>, <https://www.pydev.org/>

See <https://realpython.com/python-ides-code-editors-guide/> for a good overview of these (and other) workflow tools.



# C

# NumPy Visual Guide

**Lab Objective:** NumPy operations can be difficult to visualize, but the concepts are straightforward. This appendix provides visual demonstrations of how NumPy arrays are used with slicing syntax, stacking, broadcasting, and axis-specific operations. Though these visualizations are for 1- or 2-dimensional arrays, the concepts can be extended to  $n$ -dimensional arrays.

## Data Access

The entries of a 2-D array are the rows of the matrix (as 1-D arrays). To access a single entry, enter the row index, a comma, and the column index. Remember that indexing begins with 0.

$$A[0] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[2,1] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

## Slicing

A lone colon extracts an entire row or column from a 2-D array. The syntax  $[a:b]$  can be read as “the  $a$ th entry up to (but not including) the  $b$ th entry.” Similarly,  $[a:]$  means “the  $a$ th entry to the end” and  $[:b]$  means “everything up to (but not including) the  $b$ th entry.”

$$A[1] = A[1,:] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[:,2] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

$$A[1:,:2] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix} \quad A[1:-1,1:-1] = \begin{bmatrix} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \end{bmatrix}$$

## Stacking

`np.hstack()` stacks sequence of arrays horizontally and `np.vstack()` stacks a sequence of arrays vertically.

$$A = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

$$B = \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}$$

$$\text{np.hstack}((A, B, A)) = \begin{bmatrix} \times & \times & \times & * & * & * & \times & \times & \times \\ \times & \times & \times & * & * & * & \times & \times & \times \\ \times & \times & \times & * & * & * & \times & \times & \times \end{bmatrix}$$

$$\text{np.vstack}((A, B, A)) = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ * & * & * \\ * & * & * \\ * & * & * \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix}$$

Because 1-D arrays are flat, `np.hstack()` concatenates 1-D arrays and `np.vstack()` stacks them vertically. To make several 1-D arrays into the columns of a 2-D array, use `np.column_stack()`.

$$x = [ \times \quad \times \quad \times \quad \times ]$$

$$y = [ * \quad * \quad * \quad * ]$$

$$\text{np.hstack}((x, y, x)) = [ \times \quad \times \quad \times \quad \times \quad * \quad * \quad * \quad * \quad \times \quad \times \quad \times \quad \times ]$$

$$\text{np.vstack}((x, y, x)) = \begin{bmatrix} \times & \times & \times & \times \\ * & * & * & * \\ \times & \times & \times & \times \end{bmatrix}$$

$$\text{np.column_stack}((x, y, x)) = \begin{bmatrix} \times & * & \times \\ \times & * & \times \\ \times & * & \times \\ \times & * & \times \end{bmatrix}$$

The functions `np.concatenate()` and `np.stack()` are more general versions of `np.hstack()` and `np.vstack()`, and `np.row_stack()` is an alias for `np.vstack()`.

## Broadcasting

NumPy automatically aligns arrays for component-wise operations whenever possible. See <http://docs.scipy.org/doc/numpy/user/basics.broadcasting.html> for more in-depth examples and broadcasting rules.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} \quad x = [10 \quad 20 \quad 30]$$

$$A + x = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \\ + \\ 10 & 20 & 30 \end{bmatrix} = \begin{bmatrix} 11 & 22 & 33 \\ 11 & 22 & 33 \\ 11 & 22 & 33 \end{bmatrix}$$

$$A + x.reshape((1, -1)) = \begin{bmatrix} 1 & 2 & 3 \\ 1 & 2 & 3 \\ 1 & 2 & 3 \end{bmatrix} + \begin{bmatrix} 10 \\ 20 \\ 30 \end{bmatrix} = \begin{bmatrix} 11 & 12 & 13 \\ 21 & 22 & 23 \\ 31 & 32 & 33 \end{bmatrix}$$

## Operations along an Axis

Most array methods have an `axis` argument that allows an operation to be done along a given axis. To compute the sum of each column, use `axis=0`; to compute the sum of each row, use `axis=1`.

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{bmatrix}$$

$$A.sum(axis=0) = \left[ \begin{array}{c|c|c|c} 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \\ 1 & 2 & 3 & 4 \end{array} \right] = [4 \quad 8 \quad 12 \quad 16]$$

$$A.sum(axis=1) = \left[ \begin{array}{cccc} 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \\ \hline 1 & 2 & 3 & 4 \end{array} \right] = [10 \quad 10 \quad 10 \quad 10]$$



# D

# Matplotlib Syntax and Customization Guide

**Lab Objective:** The documentation for Matplotlib can be a little difficult to maneuver and basic information is sometimes difficult to find. This appendix condenses and demonstrates some of the more applicable and useful information on plot customizations. It is not intended to be read all at once, but rather to be used as a reference when needed. For an interative introduction to Matplotlib, see the Introduction to Matplotlib lab in Python Essentials. For more details on any specific function, refer to the Matplotlib documentation at <https://matplotlib.org/>.

## Matplotlib Interface

Matplotlib plots are made in a `Figure` object that contains one or more `Axes`, which themselves contain the graphical plotting data. Matplotlib provides two ways to create plots:

1. Call plotting functions directly from the module, such as `plt.plot()`. This will create the plot on whichever `Axes` is currently active.
2. Call plotting functions from an `Axes` object, such as `ax.plot()`. This is particularly useful for complicated plots and for animations.

Table D.1 contains a summary of functions that are used for managing `Figure` and `Axes` objects.

Function	Description
<code>add_subplot()</code>	Add a single subplot to the current figure
<code>axes()</code>	Add an axes to the current figure
<code>clf()</code>	Clear the current figure
<code>figure()</code>	Create a new figure or grab an existing figure
<code>gca()</code>	Get the current axes
<code>gcf()</code>	Get the current figure
<code>subplot()</code>	Add a single subplot to the current figure
<code>subplots()</code>	Create a figure and add several subplots to it

Table D.1: Basic functions for managing plots.

`Axes` objects are usually managed through the functions `plt.subplot()` and `plt.subplots()`. The function `subplot()` is used as `plt.subplot(nrows, ncols, plot_number)`. Note that if the inputs for `plt.subplot()` are all integers, the commas between the entries can be omitted. For example, `plt.subplot(3,2,2)` can be shortened to `plt.subplot(322)`.

The function `subplots()` is used as `plt.subplots(nrows, ncols)`, and returns a `Figure` object and an array of `Axes`. This array has the shape `(nrows, ncols)`, and can be accessed as any other array. Figure D.1 demonstrates the layout and indexing of subplots.

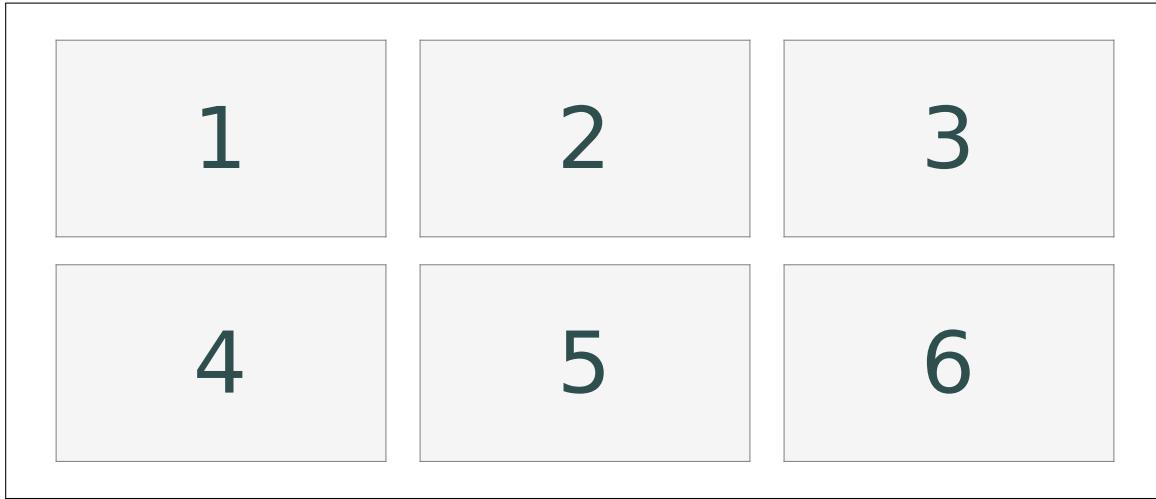


Figure D.1: The layout of subplots with `plt.subplot(2,3,i)` (2 rows, 3 columns), where `i` is the index pictured above. The outer border is the figure that the axes belong to.

The following example demonstrates three equivalent ways of producing a figure with two subplots, arranged next to each other in one row:

```
>>> x = np.linspace(-5, 5, 100)

# 1. Use plt.subplot() to switch the current axes.
>>> plt.subplot(121)
>>> plt.plot(x, 2*x)
>>> plt.subplot(122)
>>> plt.plot(x, x**2)

# 2. Use plt.subplot() to explicitly grab the two subplot axes.
>>> ax1 = plt.subplot(121)
>>> ax1.plot(x, 2*x)
>>> ax2 = plt.subplot(122)
>>> ax2.plot(x, x**2)

# 3. Use plt.subplots() to get the figure and all subplots simultaneously.
>>> fig, axes = plt.subplots(1, 2)
>>> axes[0].plot(x, 2*x)
>>> axes[1].plot(x, x**2)
```

## Achtung!

Be careful not to mix up the following similarly-named functions:

1. `plt.axes()` creates a new place to draw on the figure, while `plt.axis()` or `ax.axis()` sets properties of the *x*- and *y*-axis in the current axes, such as the *x* and *y* limits.
2. `plt.subplot()` (singular) returns a single subplot belonging to the current figure, while `plt.subplots()` (plural) creates a new figure and adds a collection of subplots to it.

# Plot Customization

## Styles

Matplotlib has a number of built-in styles that can be used to set the default appearance of plots. These can be used via the function `plt.style.use()`; for instance, `plt.style.use("seaborn")` will have Matplotlib use the "seaborn" style for all plots created afterwards. A list of built-in styles can be found at [https://matplotlib.org/stable/gallery/style\\_sheets/style\\_sheets\\_reference.html](https://matplotlib.org/stable/gallery/style_sheets/style_sheets_reference.html).

The style can also be changed only temporarily using `plt.style.context()` along with a `with` block:

```
with plt.style.context('dark_background'):
    # Any plots created here use the new style
    plt.subplot(1,2,1)
    plt.plot(x, y)
    #
# Plots created here are unaffected
plt.subplot(1,2,2)
plt.plot(x, y)
```

## Plot layout

### Axis properties

Table D.2 gives an overview of some of the functions that may be used to configure the axes of a plot.

The functions `xlim()`, `ylim()`, and `axis()` are used to set one or both of the *x* and *y* ranges of the plot. `xlim()` and `ylim()` each accept two arguments, the lower and upper bounds, or a single list of those two numbers. `axis()` accepts a single list consisting, in order, of `xmin`, `xmax`, `ymin`, `ymax`. Passing `None` instead of one of the numbers to any of these functions will make it not change the corresponding value from what it was. Each of these functions can also be called without any arguments, in which case it will return the current bounds. Note that `axis()` can also be called directly on an `Axes` object, while `xlim()` and `ylim()` cannot.

`axis()` also can be called with a string as its argument, which has several options. The most common is `axis('equal')`, which makes the scale of the *x*- and *y*-scales equal (i.e. makes circles circular).

Function	Description
<code>axis()</code>	set the $x$ - and $y$ -limits of the plot
<code>grid()</code>	add gridlines
<code>xlim()</code>	set the limits of the $x$ -axis
<code>ylim()</code>	set the limits of the $y$ -axis
<code>xticks()</code>	set the location of the tick marks on the $x$ -axis
<code>yticks()</code>	set the location of the tick marks on the $y$ -axis
<code>xscale()</code>	set the scale type to use on the $x$ -axis
<code>yscale()</code>	set the scale type to use on the $y$ -axis
<code>ax.spines[side].set_position()</code>	set the location of the given spine
<code>ax.spines[side].set_color()</code>	set the color of the given spine
<code>ax.spines[side].set_visible()</code>	set whether a spine is visible

Table D.2: Some functions for changing axis properties. `ax` is an `Axes` object.

To use a logarithmic scale on an axis, the functions `xscale("log")` and `yscale("log")` can be used.

The functions `xticks()` and `yticks()` accept a list of tick positions, which the ticks on the corresponding axis are set to. Generally, this works the best when used with `np.linspace()`. This function also optionally accepts a second argument of a list of labels for the ticks. If called with no arguments, the function returns a list of the current tick positions and labels instead.

The spines of a Matplotlib plot are the black border lines around the plot, with the left and bottom ones also being used as the axis lines. To access the spines of a plot, call `ax.spines[side]`, where `ax` is an `Axes` object and `side` is `'top'`, `'bottom'`, `'left'`, or `'right'`. Then, functions can be called on the `Spine` object to configure it.

The function `spine.set_position()` has several ways to specify the position. The two simplest are with the arguments `'center'` and `'zero'`, which place the spine in the center of the subplot or at an  $x$ - or  $y$ -coordinate of zero, respectively. The others are passed as a tuple `(position_type, amount)`:

- `'data'`: place the spine at an  $x$ - or  $y$ -coordinate equal to `amount`.
- `'axes'`: place the spine at the specified `Axes` coordinate, where 0 corresponds to the bottom or left of the subplot, and 1 corresponds to the top or right edge of the subplot.
- `'outward'`: places the spine `amount` pixels outward from the edge of the plot area. A negative value can be used to move it inwards instead.

`spine.set_color()` accepts any of the color formats Matplotlib supports. Alternately, using `set_color('none')` will make the spine not be visible. `spine.set_visible()` can also be used for this purpose.

The following example adjusts the ticks and spine positions to improve the readability of a plot of  $\sin(x)$ . The result is shown in Figure D.2.

```
>>> x = np.linspace(0,2*np.pi,150)
>>> plt.plot(x, np.sin(x))
>>> plt.title(r"$y=\sin(x)$")

#Set the ticks to multiples of pi/2, make nice labels
>>> ticks = np.pi / 2 * np.array([0,1,2,3,4])
```

```

>>> tick_labels = ["$0$", r"$\frac{\pi}{2}$", r"$\pi$", r"$\frac{3\pi}{2}$",
...                 r"$2\pi$"]
>>> plt.xticks(ticks, tick_labels)

#Move the bottom spine to zero, remove the top and right ones
>>> ax = plt.gca()
>>> ax.spines['bottom'].set_position('zero')
>>> ax.spines['right'].set_color('none')
>>> ax.spines['top'].set_color('none')

>>> plt.show()

```

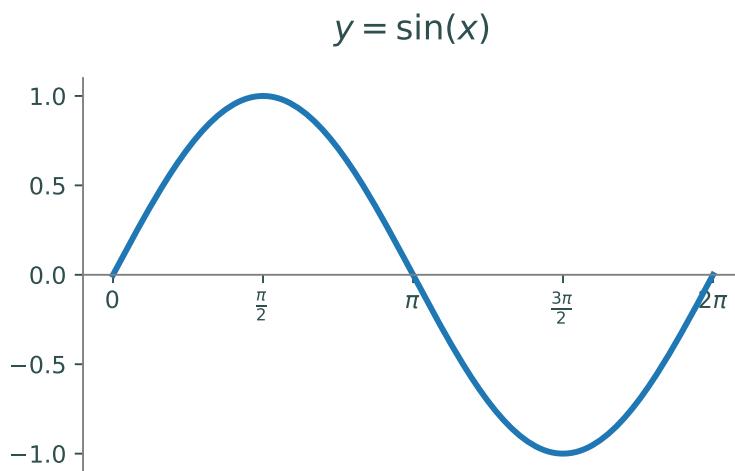


Figure D.2: Plot of  $y = \sin(x)$  with axes modified for clarity

### Plot Layout

The position and spacing of all subplots within a figure can be modified using the function `plt.subplots_adjust()`. This function accepts up to six keyword arguments that change different aspects of the spacing. `left`, `right`, `top`, and `bottom` are used to adjust the rectangle around all of the subplots. In the coordinates used, 0 corresponds to the bottom or left edge of the figure, and 1 corresponds to the top or right edge of the figure. `hspace` and `wspace` set the vertical and horizontal spacing, respectively, between subplots. The units for these are in fractions of the average height and width of all subplots in the figure. If more fine control is desired, the position of individual `Axes` objects can also be changed using `ax.get_position()` and `ax.set_position()`.

The size of the figure can be configured using the `figsize` argument when creating a figure:

```
>>> plt.figure(figsize=(12,8))
```

Note that many environments will scale the figure to fill the available space. Even so, changing the figure size can still be used to change the aspect ratio as well as the relative size of plot elements.

The following example uses `subplots_adjust()` to create space for a legend outside of the plotting space. The result is shown in Figure D.3.

```
#Generate data
>>> x1 = np.random.normal(-1, 1.0, size=60)
>>> y1 = np.random.normal(-1, 1.5, size=60)
>>> x2 = np.random.normal(2.0, 1.0, size=60)
>>> y2 = np.random.normal(-1.5, 1.5, size=60)
>>> x3 = np.random.normal(0.5, 1.5, size=60)
>>> y3 = np.random.normal(2.5, 1.5, size=60)

#Make the figure wider
>>> fig = plt.figure(figsize=(5,3))

#Plot the data
>>> plt.plot(x1, y1, 'r.', label="Dataset 1")
>>> plt.plot(x2, y2, 'g.', label="Dataset 2")
>>> plt.plot(x3, y3, 'b.', label="Dataset 3")

#Create a legend to the left of the plot
>>> lspace = 0.35
>>> plt.subplots_adjust(left=lspace)
#Put the legend at the left edge of the figure
>>> plt.legend(loc=(-lspace/(1-lspace),0.6))
>>> plt.show()
```

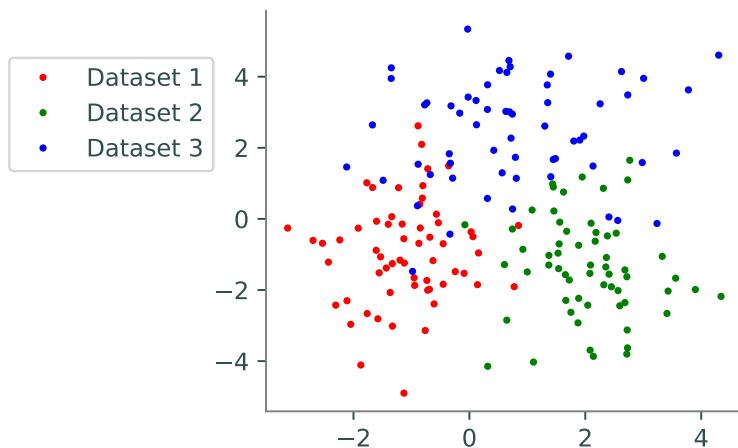


Figure D.3: Example of repositioning axes.

## Colors

The color that a plotting function uses is specified by either the `c` or `color` keyword arguments; for most functions, these can be used interchangeably. There are many ways to specify colors. The most simple is to use one of the basic colors, listed in Table D.3. Colors can also be specified using an RGB tuple such as `(0.0, 0.4, 1.0)`, a hex string such as `"#0000FF"`, or a CSS color name like `"DarkOliveGreen"` or `"FireBrick"`. A full list of named colors that Matplotlib supports can be found at [https://matplotlib.org/stable/gallery/color/named\\_colors.html](https://matplotlib.org/stable/gallery/color/named_colors.html). If no color is specified for a plot, Matplotlib automatically assigns it one from the default color cycle.

Code	Color	Code	Color
'b'	blue	'y'	yellow
'g'	green	'k'	black
'r'	red	'w'	white
'c'	cyan	'CO' - 'C9'	Default colors
'm'	magenta		

Table D.3: Basic colors available in Matplotlib

Plotting functions also accept an `alpha` keyword argument, which can be used to set the transparency. A value of 1.0 corresponds to fully opaque, and 0.0 corresponds to fully transparent.

The following example demonstrates different ways of specifying colors:

```
#Using a basic color
>>> plt.plot(x, y, 'r')
#Using a hexadecimal string
>>> plt.plot(x, y, color='FF0080')
#Using an RGB tuple
>>> plt.plot(x, y, color=(1, 0.5, 0))
#Using a named color
>>> plt.plot(x, y, color='navy')
```

## Colormaps

Certain plotting functions, such as heatmaps and contour plots, accept a colormap rather than a single color. A full list of colormaps available in Matplotlib can be found at [https://matplotlib.org/stable/gallery/color/colormap\\_reference.html](https://matplotlib.org/stable/gallery/color/colormap_reference.html). Some of the more commonly used ones are `"viridis"`, `"magma"`, and `"coolwarm"`. A colorbar can be added by calling `plt.colorbar()` after creating the plot.

Sometimes, using a logarithmic scale for the coloring is more informative. To do this, pass a `matplotlib.colors.LogNorm` object as the `norm` keyword argument:

```
# Create a heatmap with logarithmic color scaling
>>> from matplotlib.colors import LogNorm
>>> plt.pcolormesh(X, Y, Z, cmap='viridis', norm=LogNorm())
```

Function	Description	Usage
<code>annotate()</code>	adds a commentary at a given point on the plot	<code>annotate('text',(x,y))</code>
<code>arrow()</code>	draws an arrow from a given point on the plot	<code>arrow(x,y,dx,dy)</code>
<code>colorbar()</code>	Create a colorbar	<code>colorbar()</code>
<code>legend()</code>	Place a legend in the plot	<code>legend(loc='best')</code>
<code>text()</code>	Add text at a given position on the plot	<code>text(x,y,'text')</code>
<code>title()</code>	Add a title to the plot	<code>title('text')</code>
<code>suptitle()</code>	Add a title to the figure	<code>suptitle('text')</code>
<code>xlabel()</code>	Add a label to the $x$ -axis	<code>xlabel('text')</code>
<code>ylabel()</code>	Add a label to the $y$ -axis	<code>ylabel('text')</code>

Table D.4: Text and annotation functions in Matplotlib

## Text and Annotations

Matplotlib has several ways to add text and other annotations to a plot, some of which are listed in Table D.4. The color and size of the text in most of these functions can be adjusted with the `color` and `fontsize` keyword arguments.

Matplotlib also supports formatting text with L<sup>A</sup>T<sub>E</sub>X, a system for creating technical documents.<sup>1</sup> To do so, use an `r` before the string quotation mark and surround the text with dollar signs. This is particularly useful when the text contains a mathematical expression. For example, the following line of code will make the title of the plot be  $\frac{1}{2} \sin(x^2)$ :

```
>>> plt.title(r"\frac{1}{2}\sin(x^2)")
```

The function `legend()` can be used to add a legend to a plot. Its optional `loc` keyword argument specifies where to place the legend within the subplot. It defaults to `'best'`, which will cause Matplotlib to place it in whichever location overlaps with the fewest drawn objects. The other locations this function accepts are `'upper right'`, `'upper left'`, `'lower left'`, `'lower right'`, `'center left'`, `'center right'`, `'lower center'`, `'upper center'`, and `'center'`. Alternately, a tuple of  $(x,y)$  can be passed as this argument, and the bottom-left corner of the legend will be placed at that location. The point  $(0,0)$  corresponds to the bottom-left of the current subplot, and  $(1,1)$  corresponds to the top-right. This can be used to place the legend outside of the subplot, although care should be taken that it does not go outside the figure, which may require manually repositioning the subplots.

The labels the legend uses for each curve or scatterplot are specified with the `label` keyword argument when plotting the object. Note that `legend()` can also be called with non-keyword arguments to set the labels, although it is less confusing to set them when plotting.

The following example demonstrates creating a legend:

```
>>> x = np.linspace(0,2*np.pi,250)

# Plot sin(x), cos(x), and -sin(x)
# The label argument will be used as its label in the legend.
>>> plt.plot(x, np.sin(x), 'r', label=r'\sin(x)')
>>> plt.plot(x, np.cos(x), 'g', label=r'\cos(x)')
>>> plt.plot(x, -np.sin(x), 'b', label=r'-\sin(x)')
```

<sup>1</sup>See <http://www.latex-project.org/> for more information.

```
# Create the legend
>>> plt.legend()
```

## Line and marker styles

Matplotlib supports a large number of line and marker styles for line and scatter plots, which are listed in Table D.5.

character	description	character	description
-	solid line style	3	tri_left marker
--	dashed line style	4	tri_right marker
-.	dash-dot line style	s	square marker
:	dotted line style	p	pentagon marker
.	point marker	*	star marker
,	pixel marker	h	hexagon1 marker
o	circle marker	H	hexagon2 marker
v	triangle_down marker	+	plus marker
^	triangle_up marker	x	x marker
<	triangle_left marker	D	diamond marker
>	triangle_right marker	d	thin_diamond marker
1	tri_down marker		vline marker
2	tri_up marker	_	hline marker

Table D.5: Available line and marker styles in Matplotlib.

The function `plot()` has several ways to specify this argument; the simplest is to pass it as the third positional argument. The `marker` and `linestyle` keyword arguments can also be used. The size of these can be modified using `markersize` and `linewidth`. Note that by specifying a marker style but no line style, `plot()` can be used to make a scatter plot. It is also possible to use both a marker style and a line style. To set the marker using `scatter()`, use the `marker` keyword argument, with `s` being used to change the size.

The following code demonstrates specifying marker and line styles. The results are shown in Figure D.4.

```
#Use dashed lines:
>>> plt.plot(x, y, '--')
#Use only dots:
>>> plt.plot(x, y, '.')
#Use dots with a normal line:
>>> plt.plot(x, y, '.-')
#scatter() uses the marker keyword:
>>> plt.scatter(x, y, marker='+')

#With plot(), the color to use can also be specified in the same string.
#Order usually doesn't matter.
#Use red dots:
>>> plt.plot(x, y, '.r')
```

```
#Equivalent:  
>>> plt.plot(x, y, 'r.')  
  
#To change the size:  
>>> plt.plot(x, y, 'v-', linewidth=1, markersize=15)  
>>> plt.scatter(x, y, marker='+', s=12)
```

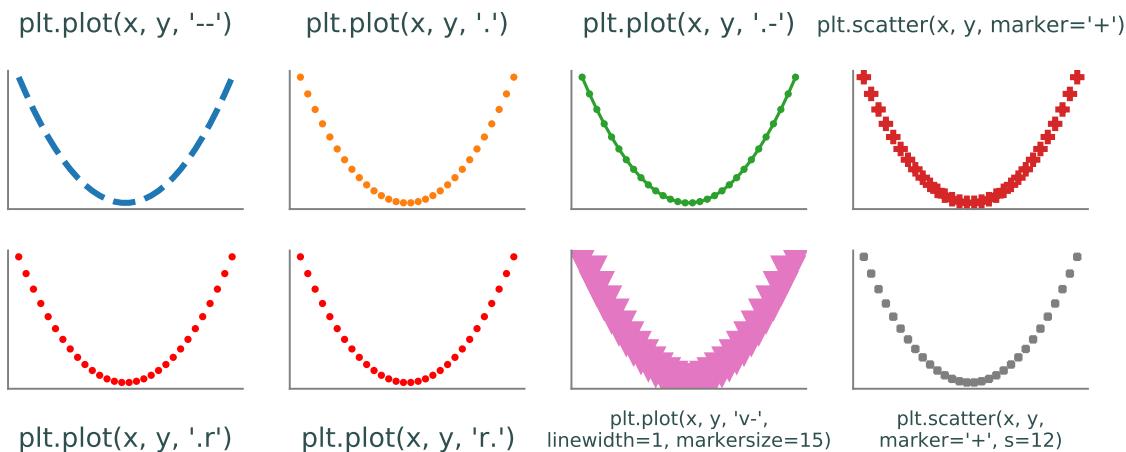


Figure D.4: Examples of setting line and marker styles.

## Plot Types

Matplotlib has functions for creating many different types of plots, many of which are listed in Table D.6. This section gives details on using certain groups of these functions.

Function	Description	Usage
<code>bar</code>	makes a bar graph	<code>bar(x,height)</code>
<code>barh</code>	makes a horizontal bar graph	<code>barh(y,width)</code>
<code>boxplots</code>	makes one or more boxplots	<code>boxplots(data)</code>
<code>contour</code>	makes a contour plot	<code>contour(X,Y,Z)</code>
<code>contourf</code>	makes a filled contour plot	<code>contourf(X,Y,Z)</code>
<code>imshow</code>	shows an image	<code>imshow(image)</code>
<code>fill</code>	plots lines with shading under the curve	<code>fill(x,y)</code>
<code>fill_between</code>	plots lines with shading between two given y values	<code>fill_between(x,y1, y2=0)</code>
<code>hexbin</code>	creates a hexbin plot	<code>hexbin(x,y)</code>
<code>hist</code>	plots a histogram from data	<code>hist(data)</code>
<code>pcolormesh</code>	makes a heatmap	<code>pcolormesh(X,Y,Z)</code>
<code>pie</code>	makes a pie chart	<code>pie(x)</code>
<code>plot</code>	plots lines and data on standard axes	<code>plot(x,y)</code>
<code>plot_surface</code>	plot a surface in 3-D space	<code>plot_surface(X,Y,Z)</code>
<code>polar</code>	plots lines and data on polar axes	<code>polar(theta,r)</code>
<code>loglog</code>	plots lines and data on logarithmic x and y axes	<code>loglog(x,y)</code>
<code>scatter</code>	plots data in a scatterplot	<code>scatter(x,y)</code>
<code>semilogx</code>	plots lines and data with a log scaled x axis	<code>semilogx(x,y)</code>
<code>semilogy</code>	plots lines and data with a log scaled y axis	<code>semilogy(x,y)</code>
<code>specgram</code>	makes a spectrogram from data	<code>specgram(x)</code>
<code>spy</code>	plots the sparsity pattern of a 2D array	<code>spy(Z)</code>
<code>triplot</code>	plots triangulation between given points	<code>triplot(x,y)</code>

Table D.6: Some basic plotting functions in Matplotlib.

## Line plots

Line plots, the most basic type of plot, are created with the `plot()` function. It accepts two lists of x- and y-values to plot, and optionally a third argument of a string of any combination of the color, line style, and marker style. Note that this method only works with the single-character color codes; to use other colors, use the `color` argument. By specifying only a marker style, this function can also be used to create scatterplots.

There are a number of functions that do essentially the same thing as `plot()` but also change the axis scaling, including `loglog()`, `semilogx()`, `semilogy()`, and `polar`. Each of these functions is used in the same manner as `plot()`, and has identical syntax.

## Bar Plots

Bar plots are a way to graph categorical data in an effective way. They are made using the `bar()` function. The most important arguments are the first two that provide the data, `x` and `height`. The first argument is a list of values for each bar, either categorical or numerical; the second argument is a list of numerical values corresponding to the height of each bar. There are other parameters that may be included as well. The `width` argument adjusts the bar widths; this can be done by choosing a single value for all of the bars, or an array to give each bar a unique width. Further, the argument `bottom` allows one to specify where each bar begins on the y-axis. Lastly, the `align` argument can be set to 'center' or 'edge' to align as desired on the x-axis. As with all plots, you can use the `color` keyword to specify any color of your choice. If you desire to make a horizontal bar graph, the syntax follows similarly using the function `barh()`, but with argument names `y`, `width`, `height` and `align`.

## Box Plots

A box plot is a way to visualize some simple statistics of a dataset. It plots the minimum, maximum, and median along with the first and third quartiles of the data. This is done by using `boxplot()` with an array of data as the argument. Matplotlib allows you to enter either a one dimensional array for a single box plot, or a 2-dimensional array where it will plot a box plot for each column of the data in the array. Box plots default to having a vertical orientation but can be easily laid out horizontally by setting `vert=False`.

## Scatter and hexbin plots

Scatterplots can be created using either `plot()` or `scatter()`. Generally, it is simpler to use `plot()`, although there are some cases where `scatter()` is better. In particular, `scatter()` allows changing the color and size of individual points within a single call to the function. This is done by passing a list of colors or sizes to the `c` or `s` arguments, respectively.

Hexbin plots are an alternative to scatterplots that show the concentration of data in regions rather than the individual points. They can be created with the function `hexbin()`. Like `plot()` and `scatter()`, this function accepts two lists of x- and y-coordinates.

## Heatmaps and contour plots

Heatmaps and contour plots are used to visualize 3-D surfaces and complex-valued functions on a flat space. Heatmaps are created using the `pcolormesh()` function. Contour plots are created using `contour()` or `contourf()`, with the latter creating a filled contour plot.

Each of these functions accepts the x-, y-, and z-coordinates as a mesh grid, or 2-D array. To create these, use the function `np.meshgrid()`:

```
>>> x = np.linspace(0,1,100)
>>> y = np.linspace(0,1,80)
>>> X, Y = np.meshgrid(x, y)
```

The z-coordinate can then be computed using the x and y mesh grids.

Note that each of these functions can accept a colormap, using the `cmap` parameter. These plots are sometimes more informative with a logarithmic color scale, which can be used by passing a `matplotlib.colors.LogNorm` object in the `norm` parameter of these functions.

With `pcolormesh()`, it is also necessary to pass `shading='auto'` or `shading='nearest'` to avoid a deprecation error.

The following example demonstrates creating heatmaps and contour plots, using a graph of  $z = (x^2 + y) \sin(y)$ . The results is shown in Figure D.5

```
>>> from matplotlib.colors import LogNorm

>>> x = np.linspace(-3,3,100)
>>> y = np.linspace(-3,3,100)
>>> X, Y = np.meshgrid(x, y)
>>> Z = (X**2+Y)*np.sin(Y)

#Heatmap
>>> plt.subplot(1,3,1)
```

```

>>> plt.pcolormesh(X, Y, Z, cmap='viridis', shading='nearest')
>>> plt.title("Heatmap")

#Contour
>>> plt.subplot(1,3,2)
>>> plt.contour(X, Y, Z, cmap='magma')
>>> plt.title("Contour plot")

#Filled contour
>>> plt.subplot(1,3,3)
>>> plt.contourf(X, Y, Z, cmap='coolwarm')
>>> plt.title("Filled contour plot")
>>> plt.colorbar()

>>> plt.show()

```

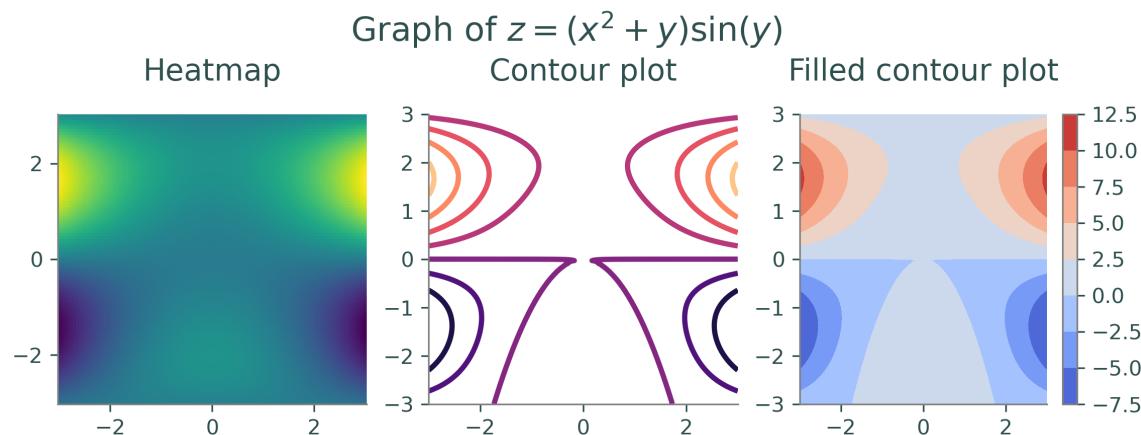


Figure D.5: Example of heatmaps and contour plots.

## Showing images

The function `imshow()` is used for showing an image in a plot, and can be used on either grayscale or color images. This function accepts a 2-D  $n \times m$  array for a grayscale image, or a 3-D  $n \times m \times 3$  array for a color image. If using a grayscale image, you also need to specify `cmap='gray'`, or it will be colored incorrectly.

It is best to also use `axis('equal')` alongside `imshow()`, or the image will most likely be stretched. This function also works best if the images values are in the range [0, 1]. Some ways to load images will format their values as integers from 0 to 255, in which case the values in the image array should be scaled before using `imshow()`.

## 3-D Plotting

Matplotlib can be used to plot curves and surfaces in 3-D space. In order to use 3-D plotting, you need to run the following line:

```
>>> from mpl_toolkits.plot3d import Axes3D
```

The argument `projection='3d'` also must be specified when creating the subplot for the 3-D object:

```
>>> plt.subplot(1,1,1, projection='3d')
```

Curves can be plotted in 3-D space using `plot()`, by passing in three lists of x-, y-, and z-coordinates. Surfaces can be plotted using `ax.plot_surface()`. This function can be used similar to creating contour plots and heatmaps, by obtaining meshes of x- and y- coordinates from `np.meshgrid()` and using those to produce the z-axis. More generally, any three 2-D arrays of meshes corresponding to x-, y-, and z-coordinates can be used. Note that it is necessary to call this function from an Axes object.

The following example demonstrates creating 3-D plots. The results are shown in Figure D.6.

```
#Create a plot of a parametric curve
ax = plt.subplot(1,3,1, projection='3d')
t = np.linspace(0, 4*np.pi, 160)
x = np.cos(t)
y = np.sin(t)
z = t / np.pi
plt.plot(x, y, z, color='b')
plt.title("Helix curve")

#Create a surface plot from np.meshgrid
ax = plt.subplot(1,3,2, projection='3d')
x = np.linspace(-1,1,80)
y = np.linspace(-1,1,80)
X, Y = np.meshgrid(x, y)
Z = X**2 - Y**2
ax.plot_surface(X, Y, Z, color='g')
plt.title(r"Hyperboloid")

#Create a surface plot less directly
ax = plt.subplot(1,3,3, projection='3d')
theta = np.linspace(-np.pi,np.pi,80)
rho = np.linspace(-np.pi/2,np.pi/2,40)
Theta, Rho = np.meshgrid(theta, rho)
X = np.cos(Theta) * np.cos(Rho)
Y = np.sin(Theta) * np.cos(Rho)
Z = np.sin(Rho)
ax.plot_surface(X, Y, Z, color='r')
plt.title(r"Sphere")

plt.show()
```

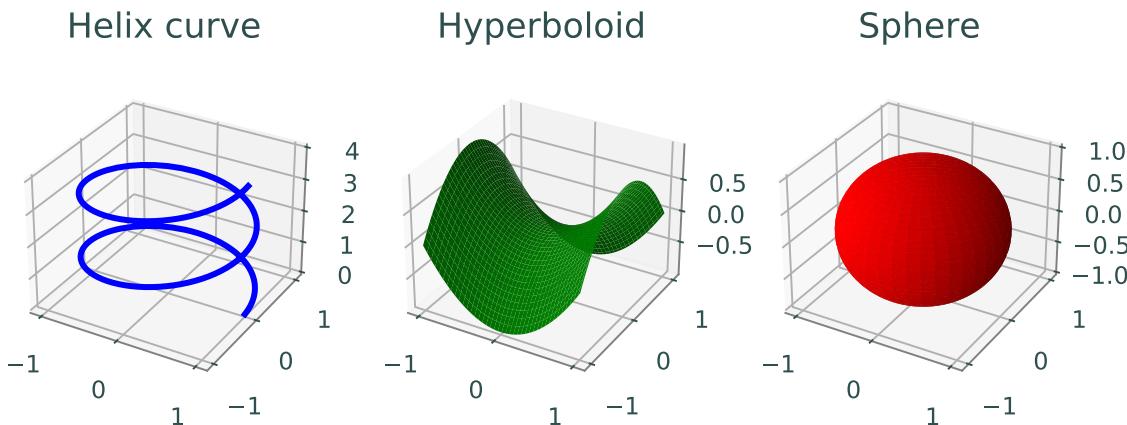


Figure D.6: Examples of 3-D plotting.

## Additional Resources

### rcParams

The default plotting parameters of Matplotlib can be set individually and with more fine control than styles by using `rcParams`. `rcParams` is a dictionary that can be accessed as either `plt.rcParams` or `matplotlib.rcParams`.

For instance, the resolution of plots can be changed via the "`figure.dpi`" parameter:

```
>>> plt.rcParams["figure.dpi"] = 600
```

A list of parameters that can set via `rcParams` can be found at [https://matplotlib.org/stable/api/matplotlib\\_configuration\\_api.html#matplotlib.RcParams](https://matplotlib.org/stable/api/matplotlib_configuration_api.html#matplotlib.RcParams).

### Animations

Matplotlib has capabilities for creating animated plots. The Animations lab in Volume 4 has detailed instructions on how to do so.

### Matplotlib gallery and tutorials

The Matplotlib documentation has a number of tutorials, found at <https://matplotlib.org/stable/tutorials/index.html>. It also has a large gallery of examples, found at <https://matplotlib.org/stable/gallery/index.html>. Both of these are excellent sources of additional information about ways to use and customize Matplotlib.



# Bibliography

- [ADH<sup>+</sup>01] David Ascher, Paul F Dubois, Konrad Hinsen, Jim Hugunin, Travis Oliphant, et al. Numerical python, 2001.
- [Hun07] J. D. Hunter. Matplotlib: A 2d graphics environment. Computing In Science & Engineering, 9(3):90–95, 2007.
- [Kim09] Seongjai Kim. Edge-preserving noise removal, part i: Second order anisotropic diffusion. Technical report, University of Kentucky Department of Mathematics, 2009.
- [Oli06] Travis E Oliphant. A guide to NumPy, volume 1. Trelgol Publishing USA, 2006.
- [Oli07] Travis E Oliphant. Python for scientific computing. Computing in Science & Engineering, 9(3), 2007.
- [PM88] Pietro Perona and Jitendra Malik. Scale-space and edge detection using anisotropic diffusion. Technical Report UCB/CSD-88-483, EECS Department, University of California, Berkeley, Dec 1988.
- [VD10] Guido VanRossum and Fred L Drake. The python language reference. Python software foundation Amsterdam, Netherlands, 2010.