

### 3. 780.20 Session 3

#### a. Follow-ups to Session 2

Here are some brief comments on various things encountered in Session 2:

- Problem 1 in Assignment #1 is to compare summing  $1/n$  from 1 to  $N$  to summing it from  $N$  to 1. When trying to understand what happens, keep in mind the simple example from Session 2 of  $1 + a + a + \dots$  vs.  $a + a + \dots + 1$ .
- When making a plot with gnuplot, it is a good idea to always `set timestamp`, which adds the current date and time to your plot. This lets you keep track when a plot was made and which of two plots is the later one.
- “Power laws,” which take the form of

$$z = Cw^\alpha , \tag{3.1}$$

with  $C$  and  $\alpha$  constants (and  $\alpha$  is often *not* an integer) are common in physics. We’ll see them frequently in looking at the errors of numerical calculations. Our simplest means of analysis is the log-log plot. Taking the logarithm of both sides of Eq. (3.1),

$$\log z = \log(Cw^\alpha) = \log C + \alpha \log w , \tag{3.2}$$

so if we plot  $y = \log z$  versus  $x = \log w$ , with  $C' = \log C$ , then we are plotting

$$y = C' + \alpha x , \tag{3.3}$$

which is a straight line with slope  $\alpha$  and intercept  $C'$ .

*Moral: If a quantity (such as an error) obeys a pure power law in some region, it will appear as a straight line in that region and the slope gives the power.*

Note that we could either output  $\log w$  and  $\log z$  from our code and plot these with gnuplot on a regular linear plot, *or* output  $w$  and  $z$  and `set logscale` to have gnuplot take the logarithms.

- **Gotcha #1.** When adding the spherical Bessel function call from GSL, you might have gotten an error message that C++ didn’t recognize the function. But there are two possible reasons to consider. If the error message was something like:

```
error: ‘gsl_sf_bessel_j1’ undeclared
```

then the problem is you didn’t include the appropriate *header file*. That is, you need the line:

```
#include <gsl/gsl_sf_bessel.h>
```

at the top of your program. The header file contains a prototype of the desired function; the GSL documentation lists the appropriate file name at the beginning of each major section. However, if the error message was something like:

```
undefined reference to ‘gsl_sf_bessel_j1’
```

then the problem was not during *compiling* but during *linking*, when the GSL library code is

combined with yours. You need to specify where this GSL code is, which is what the `-lgsl -lgslcblas` specification in the makefile (under `LDFLAGS`) does. (If you got this error, you probably didn't follow the instructions to copy a previous makefile!)

- To call C routines from C++, we have to use `extern "C"` when specifying function prototypes for the C routines in the C++ code:

```
extern "C" {
    <header stuff>
}
```

where `<header stuff>` for a single C function is just the usual prototype. (A function prototype is given at the top of the program or read in from a header file [.h file], and gives the name of the function, the number and types [e.g., `int` or `double`] of its variables and the type of its return value.) GSL routines automatically have this built into their header files, so we can call GSL routines from C++ simply by including the appropriate GSL header files.

## b. C++ Input/Output Formatting: Take 1

We've already encountered input and output to the screen in C++ using `cin` and `cout` in our first programs as well as output to a file. Here we'll elaborate a bit on their use. A good online introduction is Ref. [4].

The C++ system of input and output offers (at least) three advantages over that of C:

1. C++ has an unformatted mode (by default) for input and output. That is, we can use `cin` and `cout` without any specified formatting, as one has to do in C. This is useful for new users, for simple input/output, and for incremental code development (worry about the formatting later). This avoids a lot of bugs encountered using C's input/output functions, such as `scanf` and `printf`.
2. The C++ input/output operators can be "overloaded" to work with new data types defined in a program (e.g., a class). There is no way to do this in C with `printf` or `scanf`.
3. There is type checking in C++ output operations, unlike in C. In C, you can get bizarre output if the type of your variable and the format do not match [4], e.g.,  

```
printf("%d%s", "Hello", 27);
```

As we develop our computational physics codes, we'd like to control the format of the output, such as whether or not floating-point numbers are printed in scientific notation and how many digits are used. We can use "manipulators" to do this.

- The handout "Formatting with Manipulators" gives a summary with examples of the use of manipulators in C++.
- We can stick manipulators anywhere in an output "stream" between `<<`'s, but the placement can matter. They will only affect output following their appearance and in some cases only the next item in the stream.

- The most common manipulator is `endl`, which generates a carriage return. If you want to skip three lines:
 

```
cout << endl << endl << endl;
```

 will do it. You get use of `endl` when you include `iostream` with `#include <iostream>`
- To have access to most of the other manipulators, include the `iomanip` header file with the statement `#include <iomanip>` along with the other include files. If you don't do this, the use of manipulators such as `setprecision` will give a warning that it is undeclared.
- Here are some basic examples. See the handout for more options.
 

```
cout << rel_error << endl;
```

 $\implies$  prints the value with fixed decimal point and 5 or 6 digits.
 

```
cout << scientific << rel_error << endl;
```

 $\implies$  now in scientific notation.
 

```
cout << scientific << setprecision(16) << rel_error << endl;
```

 $\implies$  scientific with 16 digits.
 

```
cout << fixed << setprecision(6) << setw(8) << rel_error << endl;
```

 $\implies$  fixed point with precision 6 and the width set to 8 to get output to line up.

We'll see more examples as we go. See the handout on manipulators and Ref. [4] for more information.

Sending output to a file is the same as sending it to `cout`, except that we have to:

- Put `#include <fstream>` with the other include files (the “f” in “fstream” stands for “file”).
- Associate the output file (let's call it `my_file.out`) with the name of a “stream” that substitutes for `cout` (let's call it `my_out`) using `ofstream`:
 

```
ofstream my_out ("my_file.out");
```
- We can use the same conventions and manipulators as when using `cout`, e.g.,
 

```
my_out << "The relative error is " << scientific << rel_error << endl;
```

### c. Numerical Differentiation [1, 2]

The problem that round-off error causes for numerical derivatives is easy to see. In any numerical approximation to a derivative, one is simulating something like

$$\lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} = f'(x) . \quad (3.4)$$

But on a computer, the difference  $f(x+h) - f(x)$  *does not* keep decreasing with  $h$  as  $h$  decreases. Instead, it reaches a minimum approximately equal to the machine precision  $\epsilon_m$  (or equals zero). Thus,

$$f'_c(x) \xrightarrow{h \rightarrow 0} \frac{\epsilon_m}{h} \implies \log_{10} f'_c(x) \approx -\log_{10} h + \log_{10} \epsilon_m . \quad (3.5)$$

What does this imply for a graph on a log-log plot?

### c.1 Forward Difference

The easiest approximation to a numerical derivative is based on a simple Taylor expansion:

$$f(x_0 + h) = f(x_0) + hf'(x_0) + \frac{h^2}{2}f^{(2)}(x_0) + \dots, \quad (3.6)$$

and solving for  $f'(x_0)$  (we'll use a subscript  $c$  for a computed expression):

$$\begin{aligned} f'_c(x_0) &\approx \frac{f(x_0 + h) - f(x_0)}{h} \\ &\approx f'(x_0) + \frac{h}{2}f^{(2)}(x_0) + \dots \end{aligned} \quad (3.7)$$

This is called the “forward-difference” derivative algorithm. Unless the second derivative  $f^{(2)}(x_0)$  happens to vanish, the error is proportional to  $h$ .

We can improve the approximation by reducing  $h$ , but only to the point where the round-off error from the subtractive cancellation gets to be the same size. Thus, the optimal choice for  $h$  is when the round-off error is about equal to the approximation error. Since

$$\epsilon_{\text{round-off}} \approx \frac{\epsilon_m}{h} \quad \text{and} \quad \epsilon_{\text{approx}}^{\text{forward}} \approx \frac{f^{(2)}h}{2}, \quad (3.8)$$

the optimal  $h$  is about

$$h \approx \left( \frac{2\epsilon_m}{f^{(2)}} \right)^{1/2}. \quad (3.9)$$

How big is that? (Plug in some test cases.)

### c.2 Central Difference

The forward-difference algorithm approximates the slope at  $x_0$  by the slope of the line from  $x_0$  to  $x_0 + h$ . We could imagine that a better line to use has  $x_0$  in the middle rather than one of the endpoints. This is indeed the case. The “central-difference” derivative algorithm takes a half-step back and a half-step forward from  $x_0$ :

$$f'_c(x_0) \approx \frac{f(x_0 + h/2) - f(x_0 - h/2)}{h} \equiv D_c f(x, h), \quad (3.10)$$

where  $D_c$  stands for central difference.

If we use the Taylor series for  $f(x \pm h/2)$  in this formula, we find

$$f'_c(x_0) \approx f'(x_0) + \frac{h^2}{24}f^{(3)}(x_0) + \dots, \quad (3.11)$$

where all the odd powers of  $h$  cancel. So we get an extra power of  $h$  for no extra computational cost (we still only evaluate the function twice)! This is the same type of computational gain we get

from Simpson's rule (see below), which we'll see goes like  $1/N^4$  rather than the naive expectation of  $1/N^3$ .

We can apply the same error analysis as with the forward-difference algorithm to find the optimal  $h$ . Since

$$\epsilon_{\text{round-off}} \approx \frac{\epsilon_m}{h} \quad \text{and} \quad \epsilon_{\text{approx}}^{\text{central}} \approx \frac{f^{(3)}h^2}{24}, \quad (3.12)$$

the optimal  $h$  is about

$$h \approx \left( \frac{24\epsilon_m}{f^{(3)}} \right)^{1/3}. \quad (3.13)$$

Unless the derivatives are ill-behaved (so that  $f^{(3)}$  is unexpectedly large compared to  $f^{(2)}$ ), the central-difference algorithm will be far superior to the forward-difference algorithm. [Note: This means we must be very careful when taking derivatives numerically of measured data that has noise.] *Also, the optimal  $h$  is much larger for the central-difference algorithm!* We can do even better using Richardson extrapolation, which is discussed in the Session 4 notes.

#### d. Numerical Integration

Numerical integration (also called “numerical quadrature”) is both a topic worth knowing about, because it shows up frequently in computational physics problems, and is a good example of the interplay of approximation and round-off errors. The handout with an excerpt from Chapter 4 Integration in the Landau-Paez text [1] and Chapter 7 of the Hjorth-Jensen notes [2] (available from the 780 webpage under “Supplementary Readings”) provide good introductions. (Also see *Numerical Recipes* [3].) Here we'll just give some basic comments to get you started.

The basic idea [1] is that we approximate the integral of  $f(x)$  from  $a$  to  $b$  as a weighted sum of  $N$  values of the integral at selected points  $x_i$  in the interval  $[a, b]$ :

$$\int_a^b f(x) dx \approx \sum_{i=1}^N f(x_i) w_i, \quad (3.14)$$

where the  $w_i$  are a set of “weights”, which are particular to the integration rule. That is, different choices of the  $x_i$  and the  $w_i$  for a given  $N$  are what distinguish different rules. As  $N \rightarrow \infty$ , the sum formally converges to the value of the integral (but not on the computer!!). In general, the precision of the approximation increases with  $N$  until round-off errors set in (you'll explore this in Session 3).

The simplest integration rules, which we start with, have evenly spaced intervals. If each interval is small and the function is smooth, the function should be well approximated by a polynomial, which we know how to integrate term-by-term. Thus, within each interval, the function is approximated by the first terms in a Taylor series expansion of  $f$ . If we use more terms, we should get a better approximation, unless the Taylor series expansion is not valid. This happens when we have singularities in the interval like  $1/x$  or non-integral powers like  $x^{1/2}$  for  $x = 0$ . (There is no problem

for these outside of intervals containing  $x = 0$ .) We have to deal with these cases carefully, as we'll do later in Session 4.

The *trapezoid* rule approximates the function by a straight line in the interval while *Simpson's* rule approximates the function by a parabola (which should be more accurate). Check the references for derivations of the rule and the approximations. Here we'll just quote the results. The trapezoid rule is for  $N - 1$  intervals (i.e.,  $N$  points) of width

$$h = \frac{b - a}{N - 1} \quad \text{with} \quad x_i = a + (i - 1)h, \quad i = 1, \dots, N \quad (3.15)$$

is

$$\int_a^b f(x) dx \approx \frac{h}{2} f_1 + h f_2 + h f_3 + \dots + h f_{N-1} + \frac{h}{2} f_N, \quad (3.16)$$

while Simpson's rule is

$$\int_a^b f(x) dx \approx \frac{h}{3} f_1 + \frac{4h}{3} f_2 + \frac{2h}{3} f_3 + \frac{4h}{3} f_4 + \dots + \frac{4h}{3} f_{N-1} + \frac{h}{3} f_N, \quad (3.17)$$

where  $N$  must be odd. Note how these rules correspond to Table 4.1 in the Landau-Paez excerpt. The *local* and *global* approximation errors for each are discussed in Section 4.5 of that handout. The results are, for fixed  $a$  and  $b$ , that

$$\text{trapezoid} \implies \epsilon_{\text{approx}} \propto \frac{1}{N^2} \quad (3.18)$$

$$\text{Simpson's} \implies \epsilon_{\text{approx}} \propto \frac{1}{N^4}. \quad (3.19)$$

(The second result is puzzling at first; shouldn't it be  $1/N^3$ ? How do we gain an extra power of  $1/N$ ?) We'll test these results in Session 3.

The other integration rule we'll use is Gaussian quadrature, which does not use equally spaced intervals. Rather, the points  $x_i$  and weights  $w_i$  are chosen so that a polynomial of degree  $(2N - 1)$  times a specified function is integrated *exactly* over a corresponding interval. The specified function could be 1 with the interval  $[-1, 1]$  (Gauss-Legendre quadrature), or  $e^{-x}$  (Gauss-Laguerre) with the interval  $[0, \infty]$ , or many other choices. By appropriate scaling, the Gaussian quadrature intervals can be transformed to  $[a, b]$ . These integration rules are amazingly effective for a wide range of integrands, as we'll see by example in Session 3.

## e. Accumulation of Multiplicative Errors and Random Walks

We saw in Session 2 how errors combine when floating-point numbers are added or subtracted. What about when they are multiplied? Suppose  $z_1$  and  $z_2$  are multiplied to give  $z_3$ . How does the round-off error in  $z_3$  relate to the errors in  $z_1$  and  $z_2$ ? Write it out using  $z_c = z(1 + \epsilon)$  with  $|\epsilon| \leq \epsilon_m$  (recall that  $\epsilon_m$  is the machine precision and  $z_c$  means the computer representation of  $z$ ):

$$z_{3c} = z_3(1 + \epsilon_3)$$

$$\begin{aligned}
&= z_1(1 + \epsilon_1) \times z_2(1 + \epsilon_2) \\
&\doteq z_1 z_2 (1 + \epsilon_1 + \epsilon_2) \\
&\doteq z_3 (1 + \epsilon_1 + \epsilon_2) ,
\end{aligned} \tag{3.20}$$

where we've dropped  $\epsilon_1 \epsilon_2$  because it is much smaller than either  $\epsilon$  alone. Thus we find that:

$$\epsilon_3 \approx \epsilon_1 + \epsilon_2 , \tag{3.21}$$

so when we multiply (or divide, you check!), we add the round-off errors.

What is the implication of adding many round-off errors (for example, in doing a numerical integration)? It is often (usually?) the case that the errors are *uncorrelated*; that is,  $\epsilon_1$  might be anywhere in the interval  $-\epsilon_m \leq \epsilon_1 \leq \epsilon_m$  and  $\epsilon_2$  is in the same interval but is “random” compared to  $\epsilon_1$ . Let's assume that the distribution is random (the BONUS problem in problem set #1 explores the actual distribution for a test case). What do we expect to find as a total error after  $N$  multiplications?

Let  $s_i$  be a random number between  $-1$  and  $1$ , that is,

$$-1 \leq s \leq 1 . \tag{3.22}$$

Then we can say that

$$\epsilon_i = s_i \epsilon_m . \tag{3.23}$$

Now  $s_i$  should be symmetrically distributed and the width of the distribution should be about 1 (up to factors of 2), so the mean of  $s_i$  and  $s_i^2$  are:

$$\langle s_i \rangle = 0 \quad \text{and} \quad \langle s_i^2 \rangle \approx 1 . \tag{3.24}$$

On the other hand,  $s_i$  and  $s_j$  are uncorrelated, so

$$\langle s_i s_j \rangle = 0 . \tag{3.25}$$

This means that  $N$  multiplications is like a random walk with  $N$  steps. We can find the magnitude of the total error  $\epsilon_{\text{total}}$  by considering

$$\begin{aligned}
\epsilon_{\text{total}}^2 &= \epsilon_m^2 (s_1 + s_2 + \cdots + s_N)^2 \\
&= \epsilon_m^2 (s_1^2 + s_2^2 + \cdots + s_N^2 + 2s_1 s_2 + 2s_1 s_3 + \cdots) \\
&\doteq \epsilon_m^2 (N \langle s^2 \rangle + 2N \langle s_i s_j \rangle) \\
&\doteq N \epsilon_m^2 .
\end{aligned} \tag{3.26}$$

(Do you see how the third line follows from the second, since we have effectively chosen  $N$  random numbers?) Thus, the total error goes like the square root of the number of operations:

$$\epsilon_{\text{total}} \approx \sqrt{N} \epsilon_m , \tag{3.27}$$

which is characteristic of a random walk.

Bottom line: If we're in a regime where the error is dominated by round-off error rather than the approximation error, we should see the error increase as the square root of the number of operations. This tells us, for example, that we can't indefinitely improve the result of a numerical integration by making the subintervals smaller and smaller (the number of operations with scale as the number of subintervals). If we have an approximation error that goes like a power of the number of subintervals  $N$ ,  $\epsilon_{\text{approx}} \approx \alpha/N^\beta$ , with  $\beta > 0$ , then the total approximation plus round-off error will go like

$$\epsilon_{\text{total}} \approx \frac{\alpha}{N^\beta} + \sqrt{N}\epsilon_m . \quad (3.28)$$

What is the implication for the best  $N$  to choose?

## f. Pointers to Functions

Here we'll make our first (but not last) discussion of "pointers" in C++. Pointers are often a mystery even to those who use C or C++ frequently, and are usually completely baffling (at first) to those who grew up using Fortran. We'll need to revisit the use of pointers soon to take advantage of the GSL routines. In Session 3, we introduce pointers to functions, which solve the problem of how to pass a function to a subroutine. In our particular case, we want to tell a routine that does a numerical integration what integrand it should integrate.

You should note that the original "solution" to this problem in the integration program adapted from `integ.c` in the Landau–Paez book [1] was to give a name to the function, e.g., "f", in the integration subroutines and then use the *same* global name "f" to define the function in the calling program. This is not a good solution in general. For example, what if we needed to integrate more than one function? More generally, we would like a wall between the subroutine function and the calling function (a form of "encapsulation"), and all we do is pass things through the wall, without knowing precisely what is happening on the other side.

As an analogy to pointers that is now familiar to most of us, think about web pages and URL web addresses. The URL `http://www.physics.ohio-state.edu/` is clearly different from the content of the OSU Physics homepage. If we wanted to ask someone by email to look at something (e.g., a picture) on that page, we could do it two ways:

- i) send a copy of the content of the homepage, or
- ii) send them the URL.

These two ways of passing information correspond in C++ to "passing by value" and "passing by reference". The analog of the URL is the "address" of a variable or a function. The special characters `*` and `&` are used in connections with pointers. In future sessions we'll see detailed examples of their use.

For now, let's just see the how pointers to functions are used in practice. We'll take the trapezoid function in `integ_routines.cpp` as an example. In the `integ_test.cpp` program, a function called `my_integrand` is defined:

```
float my_integrand (float x)
{
    return (exp (-x));
}
```

To call trapezoid with `my_integrand` as an argument:

```
result = trapezoid (i, lower, upper, &my_integrand);
```

Note the ampersand `&` in front of `my_integrand`. That sends to trapezoid the *address* of the `my_integrand` function (i.e., it sends the location in the computer memory).

The function trapezoid is defined as:

```
float trapezoid ( int num_pts, float x_min, float x_max,
                float (*integrand) (float x) )
```

The `*` indicates a *pointer*, which holds the address that is passed in this case. For now, just note that the function

```
float integrand (float x)
```

becomes

```
float (*integrand) (float x)
```

in the function definition. (For the experts: We don't need to specify "x" in the last argument; it is a dummy variable and can be omitted entirely.)

## g. References

- [1] R.H. Landau and M.J. Paez, *Computational Physics: Problem Solving with Computers* (Wiley-Interscience, 1997).
- [2] M. Hjorth-Jensen, *Computational Physics*. These are notes from a course offered at the University of Oslo. See the 780.20 webpage for links.
- [3] W. Press *et al.*, *Numerical Recipes in C* (Cambridge, 1992). Individual chapters are available online from <http://lib-www.lanl.gov/numerical/bookcpdf.html>. There are also versions for Fortran and C++.
- [4] "Basic Input/Output" at <http://www.cplusplus.com/doc/tutorial/basic.io.html> gives a brief introduction, "C++ Notes: I/O Manipulators" gives a summary and examples at <http://www.fredosaurus.com/notes-cpp/io/omanipulators.html> and "C++ Input/Output: Streams" at <http://courses.cs.vt.edu/~cs1044/Notes/C04.IO.pdf> will step you through all the background and details.