

STEFFEN HAUG

Numerical Methods

Assignment 1

PROBLEM I

Sphere *floating* in water, i. e. the sphere is still.

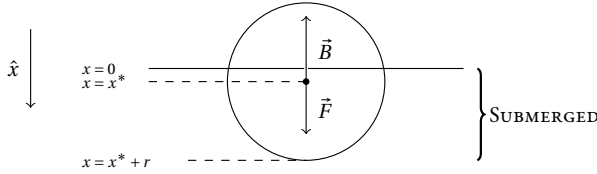


FIGURE 1: *Sphere in equilibrium.*

The sphere has a density of $\rho = 0.6$ times that of water. Archimedes' principle states that the buoyant force exerted on the sphere is equal (in magnitude) to the *weight* of the displaced fluid. Using density $\rho_{\text{fluid}} = 1$ for water, the weight of the displaced fluid is

$$-\vec{B} = m_{\text{fluid}} \cdot \vec{g} = v_{\text{fluid}} \cdot \vec{g},$$

and we know that the weight of the sphere itself is

$$\vec{F} = m_{\text{sphere}} \cdot \vec{g} = 0.6 \cdot v_{\text{sphere}} \cdot \vec{g}.$$

At equilibrium, these balance out:

$$0 = \vec{F} - \vec{B} = g(0.6v_{\text{sphere}} - v_{\text{fluid}})$$

Obviously, the volume of displaced fluid is equal to the volume of the part of the sphere that is submerged. We divide the equation by g , and thus

$$v_{\text{submerged}} = v_{\text{fluid}} = 0.6 \cdot v_{\text{sphere}} \quad (1)$$

we also know that

$$v_{\text{sphere}} = \frac{4\pi r^3}{3} \quad (2)$$

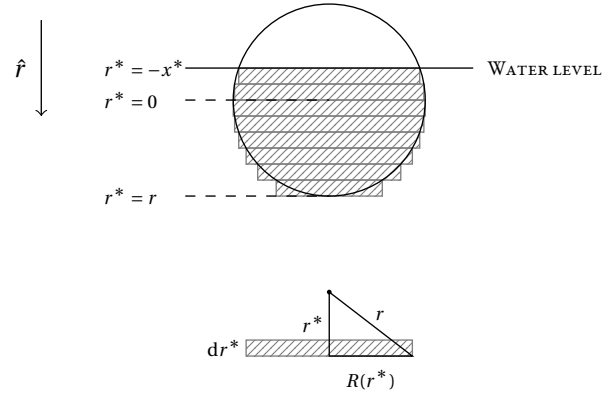


FIGURE 2: *Volume of the submerged part of the sphere.*

and with some analysis, we can work out the volume of the submerged part of the sphere. Our strategy is to integrate over cylinder-shaped volume elements.

Each volume element has a volume of $\pi \cdot R(r^*)^2 \cdot dr^*$, and we need to integrate from water level and downward in a coordinate system on the sphere. Thus,

$$v_{\text{submerged}} = \pi \int_{r^*=-x^*}^{r^*=r} R(r^*)^2 dr^* \quad (3)$$

where

$$R(r^*) = \sqrt{r^2 - (r^*)^2}.$$

We now have our mathematical model. It remains to find a solution x^* such that

$$I(x^*) := \frac{v_{\text{submerged}}}{v_{\text{sphere}}} = \frac{3}{4\pi r^3} \pi \int_{r^*=-x^*}^{r^*=r} R(r^*)^2 dr^* = \rho = 0.6$$

Given an algorithm to compute definite integrals, this can probably already be solved quite easily, but we don't know any such algorithms yet, so our strategy is to evaluate the integral by hand, and try to create a root-finding

problem.

$$\begin{aligned}
 I(x^*) &= \frac{3}{4r^3} \int_{r^*=-x^*}^{r^*=r} R(r^*)^2 dr^* \\
 &= \frac{3}{4r^3} \int_{r^*=-x^*}^{r^*=r} r^2 - (r^*)^2 dr^* \\
 &= \frac{3}{4r^3} \cdot \left[r^2 r^* - \frac{1}{3} (r^*)^3 \right]_{r^*=-x^*}^{r^*=r} \\
 &= \frac{1}{4r^3} (3r^2(r+x^*) - r^3 - (x^*)^3)
 \end{aligned}$$

This is a good time to check that

$$I(-r) = 0, \quad I(0) = 1/2, \quad \text{and} \quad I(r) = 1$$

Which is our limit cases, where the sphere is either hovering above the fluid ($x^* = -r$), or is completely submerged ($x^* = r$). The $x^* = 0$ case corresponds to exactly half the sphere being submerged.

LISTING I: *Checking for a given radius*

```

1  def I(h, r):
2      p = (3 * r**2 * (r + h) - r**3 - h**3)
3      q = (4 * r**3)
4      return p / q

>>> r = 17
>>> I(r, r)
1.0
>>> I(-r, r)
0.0
>>> I(0, r)
0.5

```

All that remains is to solve the equation

$$I(x^*) - \rho = 0; \quad 0 < \rho < 1,$$

where ρ is the density of the sphere. Our analysis reveals that $I(x^*)$ is a polynomial of degree three, which is good

because its image is all of \mathbf{R} , so a root like this is guaranteed to exist. *Unfortunately, three roots exist.* We need to define a criteria to determine which solution is the correct one.

Obviously, the sphere will not hover in mid-air, and since we have assumed $\rho < 1$, it will not sink. Clearly then, all *physical* solutions must be in the interval $[-r, r]$. If we were to prove that the CONTRACTION MAPPING THEOREM applies to some function $g: [-r, r] \rightarrow [-r, r]$, then we would have a guarantee that the physical solution is unique, and we would have defined a procedure to approximate it in the process.

APPLYING THE CONTRACTION MAPPING THEOREM

If we factor $I(x^*)$ in a particular manner, we see that

$$I(x^*) = \frac{x^*}{4r} \left(3 - \left(\frac{x^*}{r} \right)^2 \right) + \frac{1}{2},$$

and that we may write the equation $I(x^*) - \rho = 0$ as

$$x^* = g(x^*) = -4r \left(\frac{1-2\rho}{2} \right) \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-1}.$$

Notice that $g(x^*)$ is continuous on $[-r, r]$. It does have singularities at $x^* = \pm r\sqrt{3}$, but that is outside the interval, hence the mean value theorem holds:

$$g(x) - g(y) = g'(m)(x - y),$$

for some $m \in [-r, r]$. To place a bound on $g'(m)$, we need to evaluate dg/dx^* using the chain rule.

$$\begin{aligned} \frac{dg}{dx^*} &= \underbrace{-4r \left(\frac{1-2\rho}{2} \right)}_{\text{constant}} \underbrace{\frac{d}{dx^*} \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-1}}_{\gamma(x^*)} \\ \gamma &= - \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-2} \frac{d}{dx^*} \left(\frac{x^*}{r} \right)^2 \\ &= \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-2} 2 \left(\frac{x^*}{r} \right) \frac{d}{dx^*} \frac{x^*}{r} \\ &= \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-2} 2 \left(\frac{x^*}{r} \right) \frac{1}{r} \end{aligned}$$

We can rearrange this a little, and multiply by $1 = r^2 \cdot r^{-2}$:

$$\begin{aligned} \gamma &= 2x^* r^{-2} \left(3 - \left(\frac{x^*}{r} \right)^2 \right)^{-2} \cdot r^2 r^{-2} \\ &= 2x^* r^2 \left(3r^2 - \left(\frac{x^*}{r} \right)^2 r^2 \right)^{-2} \\ &= 2x^* r^2 \left(3r^2 - (x^*)^2 \right)^{-2} \end{aligned}$$

then

$$\begin{aligned} \frac{dg}{dx^*} &= -4r \left(\frac{1-2\rho}{2} \right) \gamma(x^*) \\ &= \underbrace{-8r^3 \left(\frac{1-2\rho}{2} \right)}_{\text{monotonically decreasing}} \underbrace{x^* \left(3r^2 - (x^*)^2 \right)^{-2}}_{\text{strictly positive for } x^* \in [-r, r]} \end{aligned}$$

i. e., $g'(x^*)$ is itself monotonically decreasing. This is very fortunate, since it suffices to evaluate g' in the endpoints of $[-r, r]$.

$$\begin{aligned} \left| \frac{dg}{dx^*} \right|_{x^*=-r} &= -\frac{1-2\rho}{2} \frac{8r^4}{(2r^2)^2} = 2\rho - 1 \\ \left| \frac{dg}{dx^*} \right|_{x^*=r} &= \frac{1-2\rho}{2} \frac{8r^4}{(2r^2)^2} = -(2\rho - 1) \end{aligned}$$

This finally proves that

$$\max_{x^* \in [-r, r]} \left| \frac{dg}{dx^*} \right| = |2\rho - 1|,$$

and moreover, $|2\rho - 1| < 1$ by the assumption $0 < \rho < 1$. This works out exactly so that the assumptions of the contraction mapping theorem holds:

$$\begin{aligned} |g(x) - g(y)| &= |g'(m)| |x - y| \quad \text{MEAN VALUE THEOREM} \\ &\leq |2\rho - 1| |x - y| \\ &< |x - y| \end{aligned}$$

In other words, the contraction mapping theorem holds with $L = 1$. Thus

$$\begin{cases} \{x_n\}_{n=0}^\infty \\ x_n = g(x_{n-1}) \\ x_0 \in [-r, r] \end{cases}$$

converges to the unique, physical solution, no matter the choice of x_0 .

PROBLEM 2

a. Want to use the Taylor series of f around the point a , $T_a\{f\}(x)$, to show that

$$\frac{df}{dx} = \frac{f(x + \Delta x) - f(x)}{\Delta x} + \mathcal{O}(\Delta x),$$

in other words, that the error we expect if we compute

$$\frac{df}{dx} \approx \frac{f(x + \Delta x) - f(x)}{\Delta x},$$

is proportional to Δx , i. e. the error goes to zero as Δx approaches zero.

Fix x and Δx . We expand the Taylor series to two terms, *also including the error term*, around the point x .

$$\begin{aligned} f(x^*) &= T_x \{f\}(x^*) \\ &= f(x) + f'(x)(x^* - x) + \text{error}_2(x^*) \end{aligned}$$

Which yields

$$f'(x) = \frac{f(x^*) - f(x) - \text{error}_2(x^*)}{x^* - x}$$

By evaluating the series at $x^* = x + \Delta x$

$$f'(x) = \frac{f(x + \Delta x) - f(x) - \text{error}_2(x^*)}{\Delta x}$$

To show that the error term is $\mathcal{O}(\Delta x)$ is straight forward:

$$E(\Delta x) := \frac{\text{error}_2(x^*)}{\Delta x} = \frac{1}{\Delta x} \frac{f''(\tilde{x})}{2!} (x^* - x)^2$$

for some $\tilde{x} \in [x, x^*]$. Again, we evaluate the error term at $x^* = x + \Delta x$.

$$\begin{aligned} &= \frac{f''(\tilde{x})}{2\Delta x} (\Delta x)^2 \\ &= \frac{f''(\tilde{x})}{2} \Delta x \end{aligned}$$

as long as f'' does not blow up on $I = [x, x + \Delta x]$. Because of the fact that I is bounded and closed, f'' attains its maximum and minimum on I . Thus we can let

$$L = \text{diam } f(I) = \max_{x, y \in I} |f''(x) - f''(y)|$$

Which makes

$$|E(\Delta x)| \leq L\Delta x.$$

TESTING THE APPROXIMATION

We want to compute the error committed by our formula for some function which we know how to differentiate symbolically. We want to compute

$$\left| \frac{d}{dx} \cos x \right|_{x=\pi/4}$$

We know that the symbolic derivative of $\cos x$ is $-\sin x$.

LISTING 2: *Approximate differentiation program*

```
1 # src/alg.py
2 def ddx(f, dx):
3     return lambda x: (f(x + dx) - f(x)) / dx
```

`ddx` returns a function of type `float -> float` that closes over `f` and `dx`. This new function in turn evaluates the derivative at a point `x`.

LISTING 3: *Using the differentiation program*

```
1 import src.alg as alg
2 import numpy as np
3 from math import sin, cos, pi
4
5 def errors(deltas):
6     for dx in deltas:
7         dcosdx = alg.ddx(cos, dx)
8         approx = dcosdx(pi/4)
9         exact = -sin(pi/4)
10        yield (abs(approx - exact), dx)
11
12 dxs = (1/2**k for k in range(30))
13 data = list(errors(dxs))
14 error, delta = zip(*data)
```

Plotting error against delta in a logarithmic coordinate system confirms that the error is indeed linear in Δx .

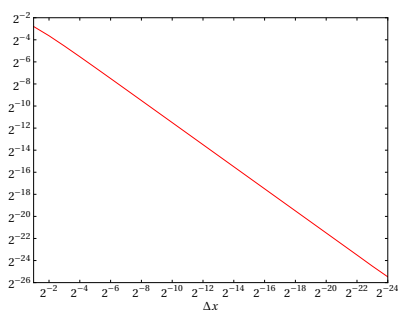


FIGURE 3: The absolute error of the approximation as a function of Δx for $f(x) = \cos x$ and $x = \pi/4$

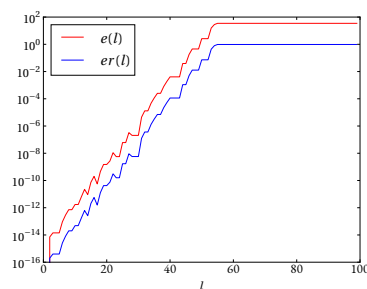


FIGURE 4: The absolute and relative error committed by the algorithm.

b. We want to observe the propagation of numerical error when iterating with a transformation using floating-point arithmetic.

LISTING 4: Iterating with sqrt

```
1 from math import sqrt
2 import random as rand
3
4 def iter_sqrt(N):
5     x = rand.random() * 100
6     xh = x
7
8     for l in range(1, N):
9         for k in range(1, l):
10             xh = sqrt(xh)
11         for k in range(1, l):
12             xh = xh**2
13
14         e = abs(x - xh)
15         er = abs(x - xh) / abs(x)
16
17         yield (e, er, l)
18
19 data = list(iter_sqrt(N = 100))
20 es, ers, ls = zip(*data)
```

PROBLEM 3

LISTING 5: Bisection and Newtons method

```
1 MAX_ITER = 100
2
3 def bisection(f, I):
4     a, b = I
5
6     assert a < b
7     assert f(a) * f(b) < 0
8
9     for _ in range(MAX_ITER):
10         m = 0.5 * (a + b)
11         yield m
12         if f(a) * f(m) < 0:
13             b = m
14         elif f(b) * f(m) < 0:
15             a = m
16         else: return # then f(m) == 0, so stop
17
18 def fpi(g, x0):
19     xn = g(x0) # x1
20     xn_m1 = x0 # x0
21     for _ in range(MAX_ITER):
22         yield xn
23         xn, xn_m1 = g(xn), xn
24
25 def newton(f, dfdx, x0):
26     g = lambda x: x - f(x) / dfdx(x)
27     yield from fpi(g, x0)
```

```

30
31 def pairs(it):
32     x = next(it)
33     y = next(it)
34     while True:
35         yield(x, y)
36         x, y = y, next(it)

```

As you can see, we can implement Newton's method as a generic fixed-point iteration. I have opted for differentiating f symbolically and providing it to `newton`.

The functions are defined in a way, such that they yield consecutively better approximations, instead of returning the best approximation in the end. This is just because we want to illustrate the convergence. Moreover, we are interested in two types of convergence, which the following figures show.

Plotting the errors versus i and j , which are the iteration count, we see clearly that Newton's method converges quadratically, while bisection converges linearly.

LISTING 6: Comparing the two methods

```

1 def g(x): return x**5 - 4*x + 2
2 def dgdx(x): return 5*x**4 - 4
3
4 n = list(alg.newton(g, dgdx, 1))
5 b = list(alg.bisect(g, (1, 5)))
6
7 def abs_err(data):
8     for xn, xn_p1 in alg.pairs(iter(data)):
9         yield abs(xn - xn_p1)
10
11 i, err_abs_n = zip(*list(enumerate(abs_err(n))))
12 _, err_abs_b = zip(*list(enumerate(abs_err(b))))
13
14 def f_err(f, data):
15     for xn in data:
16         yield abs(f(xn))
17
18 j, err_f_n = zip(*list(enumerate(f_err(g, n))))
19 _, err_f_b = zip(*list(enumerate(f_err(g, b))))

```

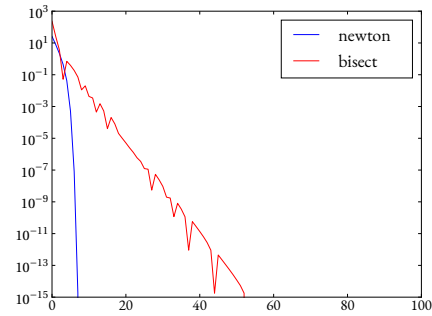


FIGURE 5: Speed of convergence $|f(x_k)|$. This clearly illustrates that the midpoint moves back and forth, as the interval shrinks from different sides.

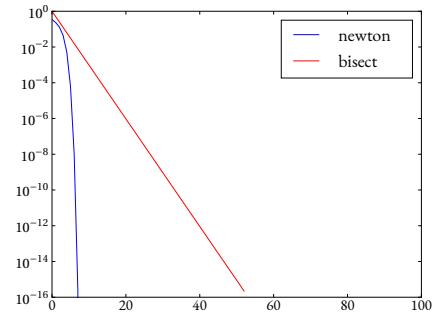


FIGURE 6: Speed of convergence $|x_{k+1} - x_k|$.

NUMERICAL EXPERIMENT

We want to verify that bisection converges linearly, and that Newton's method converges quadratically, by a numerical experiment. Recalling the definition of speed of convergence, we need to find two sequences, for bisection x_n :

$$\{\epsilon_n\}_{n=0}^{\infty}; \quad \lim_{n \rightarrow \infty} \frac{|\epsilon_{n+1}|}{|\epsilon_n|} = \mu,$$

and for Newtons method y_n :

$$\{\delta_n\}_{n=0}^\infty; \quad \lim_{n \rightarrow \infty} \frac{|\delta_{n+1}|}{|\delta_n|^2} = \nu.$$

This measure of the error is slightly different from the ones we plotted before, so some programming work is required:

LISTING 7: *Generators for the errors*

```

1 ns = list(alg.newton(g, dgdx, 1))
2 bs = list(alg.bisect(g, (1, 5)))
3
4 # define the limit as the last
5 # iteration of newtons method.
6 L = ns[-1]
7
8 ns_err = (abs(xn - L) for xn in ns)
9 bs_err = (abs(xn - L) for xn in bs)

```

BOUNDING THE ERROR

Our strategy is to try and come up with two sequences of the kind above, with the same limit, such that one bounds the error from above, and one bounds the error from below. If we can find two such sequences where the inequality holds for any number of iterations, then certainly the sequence in-between converges at the same rate. There are two problems with this: First and foremost *we only have a finite number of terms of the sequence*, so we can't say much about the asymptotic behaviour. Secondly, *we only have finite precision*, so at some point, small numbers are rounded down to zero, which means that even if we *could* compute an infinite number of elements, the asymptotic behaviour would be incorrect past some point.

We begin with the bisection method, as the linear case is (presumably) simpler. The obvious sequence to com-

pare the error with is

$$\{\epsilon_n\}(x) = \begin{cases} \epsilon_{n+1} = \mu \cdot \epsilon_n \\ \epsilon_0 = x \end{cases}$$

Where we “guess” $\mu := 1/2$, because each iteration makes the interval exactly half the size. Notice that ϵ_n can be computed with an iterative procedure using the function $x \mapsto \mu x$.

LISTING 8: *Computing the sequences*

```

1 mu = 0.5
2 es1 = list(alg.fpi(lambda x: mu * x, 10))
3 es2 = list(alg.fpi(lambda x: mu * x, 10E-4))

```

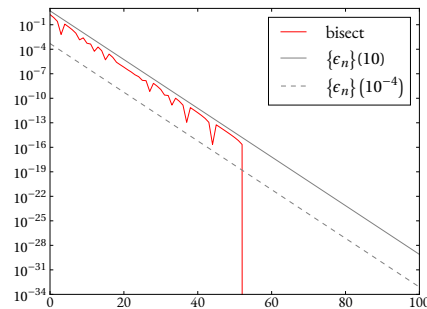


FIGURE 7: *Absolute error of bisection method, with the sequence ϵ_n . Past roughly 10^{-15} , the error “drops” suddenly.*

The sudden drop in error may be explained by the fact that the bisection algorithm needs to evaluate $f(x)$ in the endpoints, and as we approach the solution, this value becomes incredibly small. If we come *close enough*, this must be rounded down to zero, because python does not do arbitrary precision arithmetic.

For Newton’s method, the choice of constants is less clear. Never the less, we are looking for a sequence

$$\{\delta_n\}(x) = \begin{cases} \delta_{n+1} = \nu \cdot \delta_n^2 \\ \delta_0 = x \end{cases}$$

I have set $\nu := 9/10$, which i found by trial and error. Presumably, it is possible to find a better constant by doing some analysis.

We iterate with $x \mapsto \nu \cdot x^2$, and we see the same thing happen here: Errors smaller than 10^{-15} or so are rounded down to zero, so it is really hard to see if the error is *really* bounded by the sequences. Not only that, but Newtons method converges so quickly that we can only consider the first seven iterations before the error is rounded down.

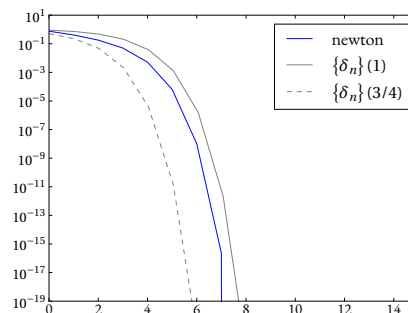


FIGURE 8: *Absolute error of newtons method, with the sequence δ_n .*