

MTH 351 Homework 3

Philip Warton

February 2, 2020

1.

(a)

We want to find a polynomial $P(x)$ such that $\max_{x \in [2,4]} |\cos(x^2) - P(x)| < 10^{-3}$. Note that $\cos(t) = 1 - \frac{t^2}{2!} + \frac{t^4}{4!} - \frac{t^6}{6!} + \dots$. We can break this into a finite Taylor polynomial and an error term as $\cos(t) = 1 - \frac{t^2}{2!} + \frac{t^4}{4!} - \frac{t^6}{6!} + \dots + (-1)^n \frac{t^{2n}}{(2n)!} + R_n(t)$ where $R_n(t)$ is our error. It should be stated that $|\cos(f(x))| \leq 1$ for any function $f(x)$ and similarly $|\sin(f(x))| \leq 1$. The Lagrange error will be in the form $\left| \frac{f^{n+1}(c)}{(n+1)!} (x - x_0)^{n+1} \right|$. Since all derivatives of cosine are sine or cosine, we can write

$$\begin{aligned} \left| \frac{f^{n+1}(c)}{(n+1)!} (t - t_0)^{n+1} \right| &\leq \frac{1}{(n+1)!} |(t - t_0)^{n+1}| \\ &= \frac{1}{(n+1)!} |(t - t_0)^{n+1}| \quad \text{let } t = x^2, t_0 = 0 \text{ and } x = 4 \in [2, 4] \\ &= \frac{1}{(n+1)!} |(4^2)^{n+1}| \\ &= \frac{1}{(n+1)!} |(16)^{n+1}| < \frac{1}{10^3} \end{aligned}$$

Using computer calculations, we get $n \geq 47$. The following is code used to get this value:

```
for n in range(0, 1000):
    R_bound = (4.0**((2.0*n)+2.0))/(math.factorial(n+1))
    if (R_bound < float(10**(-3))):
        print(n)
        break
```

Since our polynomial only accounts for the terms where n is even, we must take $\frac{47+1}{2}$ when we build our polynomial in the polynomial form of $\cos(t)$ where $t = x^2$. Thus we have $P(x) = \sum_{n=0}^{24} (-1)^n \frac{x^{4n}}{(2n)!}$ which meets the stated error term bounds.

(b)

2.

Let our model for numbers be as follows: $c_0 \ b_1 \ b_2 \ b_3 \ b_4 \ a_1 \ a_2 \ a_3$ where c_0 is the sign part, $b_1 \ b_2 \ b_3 \ b_4$ is the exponent part, and $a_1 \ a_2 \ a_3$ is the mantissa part.

(a)

Find the dynamic range and the machine epsilon.

To find the dynamic range let us find both the largest and smallest strictly positive numbers that can be represented by our model. First we wish to find the smallest number. Let $c_0 = 0$ so that our number is positive. Let $b_1b_2b_3b_4 = 0000$ so that our exponent will be its lowest possible value 2^{-6} . Finally let $a_1a_2a_3 = 001$ such that the number will be as small as possible without being zero. From this we have

$$\begin{aligned}x &= (1)(0.001)_2(2^{-6}) \\&= 1 * 2^{-3} * 2^{-6} \\&= 2^{-9}\end{aligned}$$

Therefore our smallest number is 2^{-9} .

To find the largest possible number let us maximize each component without x being infinity. Of course let $c_0 = 0$ as our number must be positive. Then let $b_1b_2b_3b_4 = 1110$ as this as large as our exponent can be without making x represent infinity. From there, let $a_1a_2a_3 = 111$ so that it is as large as possible. This gives us

$$\begin{aligned}x &= (1)(1.111)_2(2^7) \\&= (2 - 2^{-3})(2^7) \\&= \left(\frac{15}{8}\right)(2^7)\end{aligned}$$

Given we have $x_{\min} = 2^{-9}$ and $x_{\max} = \left(\frac{15}{8}\right)(2^7)$, we can say that the dynamic range is about $\frac{\left(\frac{15}{8}\right)(2^7)}{2^{-9}} = \left(\frac{15}{8}\right)2^{16} = 122880$

Now to find the machine epsilon we must take the smallest number larger than 1 and find the difference between 1 and that number. As our exponent increases, the precision of our number decreases. If our exponent is less than 0, then our number is guaranteed to be less than one, therefore let our exponent equal zero. To make our exponent zero, E must equal 7. Let our mantissa part be the smallest possible non-zero number. Altogether we have $c_0 = 0$, $b_1b_2b_3b_4 = 0111$, and $a_1a_2a_3 = 001$. We then get

$$\begin{aligned}x &= (1)(1.001)_2(2^0) \\&= (1.001)_2 \\&= 1 + \frac{1}{2^3} \\&= \frac{9}{8}\end{aligned}$$

Thus $1 - \frac{9}{8} = \frac{1}{8}$ is our machine epsilon.

(b)

Write the numbers represented by the following bit sequences:

11001001

From this we have $c_0 = 1$ therefore the sign will be negative. Our exponent part will be $2^{(1001)_2 - 7}$. For the mantissa part we have $(1.001)_2$. This gives us $(-1)\left(\frac{9}{8}\right)(2^{9-7}) = \frac{-9}{8} * 2^2 = \frac{-9}{2} = -4.5$.

00000000

Therefore we have a sign part as (1), and with our exponent part $E = 0$, we have our exponent part as 2^{-6} and no leading 1 on our mantissa part. Thus we have $(0.000)_2$ as our mantissa part. With $x = (1)(0.000)_2(2^{-6}) = 0$, we have this number representing 0.

11111000

With the sign part being (-1) and $(1111)_2 = E = 15$, we have that $x = -\infty$, regardless of mantissa part.

3.

(a)

Represent the following numbers in the style presented in problem 2.

$$x = 1$$

$$\begin{aligned} x &= 1 \\ &= (1)(1)(1) \\ &= (1)(1.000)_2(2^0) \end{aligned}$$

Therefore let $c_0 = 0$, $b_1b_2b_3b_4 = 0111$, $a_1a_2a_3 = 000$ and our number is 00111000.

$$x = 5.5$$

$$\begin{aligned} x &= 5.5 \\ &= (101.1)_2 \\ &= (1.011)_2(2^2) \\ &= (1)(1.011)_2(2^2) \end{aligned}$$

Thus we have $c_0 = 0$, $b_1b_2b_3b_4 = 1001$, $a_1a_2a_3 = 011$ and our number is 01001011.

$$x = 12.9$$

$$\begin{aligned} x &= 12.9 \\ &= (1100.1110011\dots)_2 \\ &= (1.10011100\dots)_2(2^3) \\ &\approx (1)(1.101)_2(2^3) \end{aligned}$$

Thus we have $c_0 = 0$, $b_1b_2b_3b_4 = 1010$, $a_1a_2a_3 = 101$ and our number is 01010101.

$$x = 1000$$

From problem 2 we know that the largest possible number we can represent is $15 * 2^4 = 240$. Therefore, we must use the exponent part representing infinity, and with the number being positive it must be a positive infinity. Given this we can write our sign part as (-1) , our exponent part as 1111, and our mantissa part as any collection of numbers. The number can be written as $11111a_1a_2a_3$ where a_n can be either 1 or 0.

$$x = 0.0001$$

We can write this number as $\frac{1}{10000}$. The smallest possible number we can represent we know to be $2^{-9} = \frac{1}{512}$. Therefore this number is too small and will be rounded to zero. We can write this number as 00000000.

(b)

Find the smallest number that is larger than the following numbers, as in the format from problem 2.

$x > 5.5$ To find the smallest number greater than 5.5, let us add one to our mantissa part, and leave the sign and exponent parts alone. If we were to change the sign, our number would be less than 5.5, and if we were to change the exponent, our number would be off by at least one binary place. Therefore we replace our former mantissa part 011 with 100. This gives us $(1.100)_2 * 2^2 = (6)_{10}$. This number is written as 01001111.

$x > 12.9$ Due to rounding, our previous representation of 12.9 was equal to 13 and therefore larger than 12.9. For our purposes, however, let us find the next number larger than that. Once again we must take our previous number and add one to our mantissa part. This gives us $(1.110)_2 * 2^3 = 14$. This number is represented by 01010110.

$x > 100.25$ Let us find the binary representation of 100.25. We can write $(100.25)_{10} = (1100100.01)_2$. Our representation of this would be $(1.101)_2 * 2^6$. This would make our exponent part (1101). The smallest number greater than this would be $(1.110)_2 * 2^6 = 104$. We write this as 01101110.