# Math 125 Hw 1

**1a)** Cancellation at $x=0$ as $f(x) = \frac{0}{0}$

To fix: $\dfrac{\sqrt{25+x}-5}{x}\left(\dfrac{\sqrt{25+x}+5}{\sqrt{25+x}+5}\right)$

$= \dfrac{x}{x\sqrt{25+x}+5} = \dfrac{1}{\sqrt{25+x}+5} = g(x)$

$\lim\limits_{x\to 0}\dfrac{1}{\sqrt{25+x}+5} = \dfrac{1}{10}$ $\qquad$ $\lim\limits_{x\to 0}\dfrac{\sqrt{25+x}-5}{x} = \lim\limits_{x\to 0}\dfrac{1}{2\sqrt{25+x}} = \dfrac{1}{10}$

**b)** Cancellation at $x=0$ $\quad f(x) = \frac{0}{0}$

To fix, note for small values of $x$, $\sin x = x$.

So, $f(x) = \dfrac{1-\cos x}{\sin^2 x} = \dfrac{1-\cos x}{1-\cos^2 x} = \dfrac{1-\cos x}{(1-\cos x)(1+\cos x)} = \dfrac{1}{1+\cos x}$

$\lim\limits_{x\to 0}\dfrac{1}{1+\cos x} = \dfrac{1}{2}$ $\qquad$ $\lim\limits_{x\to 0}\dfrac{1-\cos x}{x^2} = \lim\limits_{x\to 0}\dfrac{\sin x}{2x} = \lim\limits_{x\to 0}\dfrac{\cos x}{2} = \dfrac{1}{2}$ ✓

**c)** Catastrophic cancellation for $x = \pm 2\pi n$ where $n \in \mathbb{N}$.

$f = \frac{0}{0}$

To fix $\lim\limits_{x\to \pm 2\pi n}\dfrac{1-\sec(x)}{1-\sec^2 x} \Rightarrow \lim\limits_{x\to 0}\dfrac{1-\sec x}{(1-\sec x)(1+\sec x)} \quad \lim\limits_{x\to \pm 2\pi n}\dfrac{1}{1+\sec x} = \dfrac{1}{2}$

So use $\dfrac{1}{1+\sec x}$ which is same form, via trig

identities, and $\lim\limits_{x\to \pm 2\pi n}\dfrac{1-\sec x}{\tan^2 x} = \lim\limits_{x\to \pm 2\pi n}\dfrac{1}{1+\sec x} = \dfrac{1}{2}$

**2a)** $\det(A) = \det(L)\det(U)$

$\det(A) = (L_1 \cdot L_2 \cdots L_k)(U_1 \cdot U_2 \cdots U_k)$

Where $L_1 \ldots L_k$ and $U_1 \ldots U_k$ are diagonal element.

To fix under/overflow, take log:

$\log(\det A) = \log(L_1 \cdot L_2 \cdots L_k \cdot U_1 \cdot U_2 \cdots U_k)$

$$\log(\det A) = \sum_{i=1}^{K} \log L_i + \log U_i$$

This handles under/overflow, since it scales the numbers, like $\log 10^{-26} = -26$, and $\log 10^{26} = 26$, the original number can be easily recovered by doing $10^{\log \det A}$.

2b) Let $M = \max(x_1, x_2, \ldots x_n)$

$$\vec{X} = \vec{X} \cdot \frac{M}{M} = \begin{bmatrix} \frac{M}{M} x_1 \\ \vdots \\ \frac{M}{M} x_n \end{bmatrix}$$

So $\|\vec{X}\| = \sqrt{\left(\frac{M}{M} x_1\right)^2 + \left(\frac{M}{M} x_2\right)^2 + \ldots \left(\frac{M}{M}\right) x_n)^2}$

$$= M \sqrt{\left(\frac{x_1}{M}\right)^2 + \left(\frac{x_2}{M}\right)^2 + \ldots \left(\frac{x_n}{M}\right)^2}$$

- This handles overflow, however, underflow isn't really an issue with finding $\|\vec{X}\|$.

3a) $10|110 = 46$ in base 10

b) $0.375 = 0.011$      $.375 = \frac{1}{4} + \frac{1}{8}$
$1.25 = 1.010$
$1.25 + .375 = 1.625$.

In binary:
$$\begin{array}{r} \overset{1}{0.011} \\ + 1.010 \\ \hline 1.101 \end{array} = 1.625 = \text{base } 10$$

# M125

3c) For float: [1] . [23] [8]

sign    fraction   exponent

$$x = (-1)^s \, 2^{e-127} \, (1+f)$$

For smallest:

$e=1$, as $e \neq 00s$, so smallest $= 2^{1-127} (1+0) = 2^{-126}$

$f=0$

$= 1.175 \times 10^{-38}$

Largest:

All ones, minus the first one, so $e = 254$. as all $1s = 255$

$$\text{largest} = 2^{254-127} \cdot (1 + 2^{-1} + 2^{-2} + 2^{-3} + \dots 2^{-23})$$

$$= 2^{127} \left(1 + \sum_{i=1}^{23} \left(\tfrac{1}{2}\right)^i \right)$$

$$= 2^{127} \left(1 + \frac{1 - \left(\tfrac{1}{2}\right)^{24}}{1 - \tfrac{1}{2}} \right)$$

$$= 2^{127} \cdot 1.9999998 = 3.40 \times 10^{38}$$

Max float: $3.403 \times 10^{38}$

Min float: $1.175 \times 10^{-38}$

d) $\pi$ doesn't have an exact representation as it is irrational

e) $\forall$ float $\in [1,2]$, we multiply by $2^k$. Our domain is now $[1 \cdot 2^k, 2 \cdot 2^k] = [2^k, 2^{k+1}]$ Since we multiply all floats by $2^k$ no new floats are introduced. However, there are the same amount of floats, and as for $k > 0$, $2^{k+1} - 2^k > 1$, then the gaps between large numbers is bigger than the gaps between small numbers.

4 To show $\frac{|y-f(C_y)|}{|y|} \leq B^{-m}$, First let's

determine an upper band for $|y-f(C_y)|$

$$y = \pm B^c \sum_{k=\infty}^{} d_k B^{-k} = \pm \left(B^c \sum_{k=\infty}^{m+1} d_k B^{-k} + B^c \sum_{k=m}^{} d_k B^{-k}\right)$$

↳ Finite geometric convergence

$$B_c \sum_{k=\infty}^{m+1} d_k B^{-k} \leq B^c \sum_{k=\infty}^{m+1} (B-1) B^{-k}$$

$$\leq B^c (B-1) \sum_{k=\infty}^{m+1} B^{-k}$$

$$\leq B^c (B-1) \cdot \frac{B^{-(m+1)}}{\left(\frac{B-1}{B}\right)}$$

$$B_c \sum_{k=\infty}^{m+1} d_k B^{-k} \leq \frac{B^c (B-1)}{B} = B^c B^{-m}$$

→ Since in denominator

For an upper band on $|y|$, $|y| \geq B_c \sum_{k=0}^{\infty} d_k B^{-k} \geq B^c$

as minimally, $|y| = B^c B^{-k}$ as $d_1 \ldots d_k = 0$

So $\frac{|y-f(C_y)|}{|y|}$ Note for $|y-f(C_y)|$, $B^c \sum_{k=m}^{} d_k B^{-k}$

cancel so we just use our bands

$$\frac{|B^c B^{-m}|}{|B^c|} \leq B^{-m} \rightarrow B^{-m} \leq B^{-m} \quad ☐$$

5 a) $p(x) = x^5 (3x^2 + 2x^3 - 2x^6 + 9x^9)$

$= x^5 (x^2 (3 + 2x - 2x^4 + 9x^7))$

$= x^5 (x^2 (3 + x(2 - 2x^3 + 9x^6)))$

$= x^5 (x^2 (3 + x(2 + x^3(-2 + 9x^3))))$

$p(x) = x^3 \cdot x^2 (x^2 (3 + x (2 + x^3 (-2 + 9x^2 \cdot x))))$

Store $x^2$ and $x^3$ upon initial computations for reuse.

5b) $C_k = a_k$

$C_{k-1} = a_{k-1} + a_k X_0$

$C_{k-2} = a_{k-2} + (a_{k-1} + a_k X_0) X_0$

$C_{k-3} = a_{k-3} + (a_{k-2} + a_{k-1} X_0 + a_k X_0^2) X_0$

$\downarrow$ Following this pattern for later

$C_{k-k} = a_{k-k} + (a_{k-(k+1)} + a_{k-(k+2)} X_0 + \ldots + a_k X_0^{k-1}) X_0$

$C_0 = a_0 + (a_1 + a_2 X_0 + \ldots + a_k X_0^{k-1}) X_0$

$C_0 = a_0 + a_1 X_0 + a_2 X_0^2 + \ldots + a_k X_0^k$

Since $P$ is a polynomial, the most efficient way to factor is Horner's Scheme which takes $2k$ operations for a degree $k$ polynomial.

c) See code for specifics. My code got value of $6891205$, which matches calculator. Code at end of PDF on last page.

d) The most efficient method would be the following steps.
1.) $X+1$
2.) $(X+1)^2$
3.) $\left((X+1)^2\right)^2 = (X+1)^4$
4.) $\left((X+1)^4\right)^2 = (X+1)^8$

This doesn't contradict Horner's Scheme as Horner's Scheme is most efficient w/ no structure in $p(x)$.

6a) If backwards stable, $f(x,y) = f(\tilde{x}, \tilde{y})$

$fl(x) \oplus fl(y) = [x(1+\epsilon_1) + y(1+\epsilon_2)](1+\epsilon_3)$

Let $\tilde{x} = x(1+\epsilon_1)(1+\epsilon_3)$, $\tilde{y} = y(1+\epsilon_2)(1+\epsilon_3)$

We must prove $\frac{|\tilde{X}-x|}{|x|} \leq \epsilon_{mo}$ and $\frac{|\tilde{y}-y|}{|y|} \leq \epsilon_m$.

$$\frac{|\tilde{X}-x|}{|x|} = \frac{|x(1+\epsilon_1)(1+\epsilon_3)-x|}{|x|} = \frac{|x+x\epsilon_1+x\epsilon_3+x\epsilon_1\epsilon_3-x|}{|x|}$$

$$= \frac{|x\epsilon_1+x\epsilon_3+x\epsilon_1\epsilon_3|}{|x|} = \epsilon_1+\epsilon_3+\epsilon_1\cdot\epsilon_3$$

Since $\epsilon$ is small,

$$|\epsilon_1+\epsilon_3+\epsilon_1\epsilon_3| \leq |\epsilon_1|+|\epsilon_3|+|\epsilon_1\epsilon_3|$$
$$< 3\epsilon_{machine}$$
$$\leq O(\epsilon_{machine})$$

For $y$, it's a similar process.

$$\frac{|\tilde{y}-y|}{|y|} = \frac{|y(1+\epsilon_2)(1+\epsilon_3)-y|}{|y|} = \frac{|y\epsilon_2+y\epsilon_3+y\epsilon_2\epsilon_3|}{|y|}$$

$$= |\epsilon_2+\epsilon_3+\epsilon_2\epsilon_3| \leq |\epsilon_2|+|\epsilon_3|+|\epsilon_2\epsilon_3|$$
$$\leq 3\epsilon_{machine}$$
$$\leq O(\epsilon_{machine})$$

As a result, $f(x,y) = f(\tilde{x},\tilde{y})$ and subtraction is backwards stable.

b) $Xy^T = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} [y_1 \cdots y_n]$

$$= \begin{bmatrix} x_1y_1(1+\epsilon_1) & \cdots & x_1y_n(1+\epsilon_j) \\ & & \\ x_ny_1(1+\epsilon_n) & \cdots & x_ny_n(1+\epsilon_m) \end{bmatrix}$$

Each multiplication has it's own $\epsilon$.

Normally, $x \cdot y^T = \begin{bmatrix} x_1y_1 & \cdots & x_1y_m \\ & \vdots & \\ x_ny_1 & \cdots & x_ny_m \end{bmatrix} = $ rank 1.

However by introducing $(1+\epsilon)$ we don't preserve rank anymore, as $xty^T$ can no longer be rank 1, as each product is multiplied by a unique $(1+\epsilon)$ term.

7a) Number converted to binary is:
0.09999996463256883598

−I'm calling it 0.0999999 for simplicity

$$\left| 0.0999999 - \frac{1}{10} \right| = \boxed{0.0000001}$$

b) $\dfrac{.0000001 \text{ error}}{1/10 \text{ sec}} = \dfrac{.0000001 \text{ error}}{\text{sec}} \cdot \dfrac{60 \text{ sec}}{\text{min}} \cdot \dfrac{60 \text{ min}}{\text{hr}} \cdot \dfrac{100 \text{ hr}}{\text{operation}}$

$= 0.036 \text{ hrs off} = \boxed{129.6 \text{ seconds}}$

c) $\dfrac{3750 \text{ miles}}{1 \text{ hr}} \cdot \dfrac{1 \text{ hr}}{60 \text{ min}} \cdot \dfrac{1 \text{ minute}}{60 \text{ sec}} = 1.04 \text{ miles}$

$1.04 \text{ miles} \times 129.6 \text{ sec} = \boxed{134.784 \text{ miles}}$

```python
#This function implements Horner's method in Python
def horner(coeff, x):
    result = coeff[0]
    for i in range(1, len(coeff)):
        result = result*x+coeff[i]
    return result

#p(x) = 7x^3 -11x^2 +12x + 5 is polynomial in question
#Horner's method means p(x) = 5 + x(12 + x(-11+7x)))

coeff = [7, -11, 12, 5] #Coefficients of p(x) for use
x = 100 #Initial value
print(horner(coeff, x)) #Prints result
```