

CS370 Notes

Minyang Jiang

January 16, 2017

1 Floating Point Number Systems

1.1 Introduction

$$F(\beta, t, L, u)$$

continuas:

$$\begin{aligned} &0 \text{ or } \pm 0.\beta_1\beta_2\dots\beta_t * \beta^d \\ &\text{where } \beta_1 \neq 0 \\ &0 \leq \beta_i \leq \beta \\ &L \leq d \leq u \end{aligned}$$

Two common systems used today - Base 2

1. Single precision: $F(2, 24, -126, 127)$
2. Double precision: $F(2, 53, -1022, 1023)$

Important concepts

If x is any real number then set $fl(x)$ = floating point representation of x
if we write:

$$\begin{aligned} x &= \pm 0.x_1x_2x_3\dots x_tx_{t+1}\dots * \beta^d \\ fl(x) &= \pm 0.x_1x_2x_3\dots x_t * \beta^d \end{aligned}$$

relative error

$$\begin{aligned} \delta_x &= \frac{fl(x) - x}{x} \\ \|\delta_x\| &\leq ? \end{aligned}$$

$$\begin{aligned} \frac{\|fl(x) - x\|}{\|x\|} &= \frac{0.00\dots 0x_{t+1}\dots * \beta^d}{0.x_1x_2\dots x_{t+1}\dots * \beta^d} \\ &= \frac{0.x_{t+1}x_{t+2}\dots * \beta^{-t}}{x_1.x_2\dots * \beta_{-1}} \\ \delta &= \frac{fl(x) - x}{x} \quad \text{then } \|\delta\| \leq \epsilon = \begin{cases} \beta^{1-t} \\ \frac{\beta^{1-t}}{2} \end{cases} \\ fl(x) &= x(1 + \delta) \quad |\delta| \leq \epsilon \end{aligned}$$

What about floating point arithmetic?

x, y real numbers, $x + y$ real

$x \oplus y$ = addition inside floating pt system = $fl(fl(x) + fl(y))$

1.2 Analysis some errors in computation

Example. Addition

$$\begin{aligned}
 \left\| \frac{(x+y) - (x \oplus y)}{x+y} \right\| &= \frac{\|(x+y) - fl(fl(x) + fl(y))\|}{x+y} = \frac{\|x+y - x(1+\delta) + y(1+\delta)\|}{x+y} \\
 &= \frac{\|x+y - (x+y + \delta_1x + \delta_2y + x\delta_3 + y\delta_3 + \delta_1\delta_3x + \delta_2\delta_3y)\|}{x+y} \\
 &\leq \frac{\|\delta_1x\| + \|\delta_2y\| + \|\delta_3x\| + \|\delta_3y\| + \|\delta_1\delta_3x\| + \|\delta_2\delta_3y\|}{\|x+y\|} \\
 &\leq \frac{(\|x\| + \|y\|)(2\epsilon + \epsilon^2)}{\|x+y\|} \\
 \|\delta_1\| &\leq \epsilon, \quad \|\delta_2\| \leq \epsilon, \quad \|\delta_3\| \leq \epsilon
 \end{aligned}$$

\Rightarrow if x and y have same sign then relative error of addition

$$\left\| \frac{x \oplus y - (x+y)}{x+y} \right\| \leq 2\epsilon + \epsilon^2$$

However if x and y have opposite sign, then you potentially have a problem particularly when $x+y \approx 0$, Situation is called **Catastrophic cancellation**

$$\begin{aligned}
 x &= 0.x_1x_2 \dots x_{t-1}x_t x_{t+1} \dots * \beta^d \\
 y &= -0.x_1x_2 \dots x_{t-1}x_t x_{t+1} \dots * \beta^d \\
 x+y &= 0.00 \dots 0?? \dots * \beta^d
 \end{aligned}$$

1.3 How about some algorithms?

1.3.1 Example

Given α , compute:

$$I_n = \int_0^1 \frac{x^n}{x+\alpha} dx \quad n = 0, 1, \dots, 100 \dots$$

Step 1

$$I_0 = \int_0^1 \frac{1}{x+\alpha} dx = \ln(x+\alpha) = \ln(1+\alpha) - \ln(\alpha) = \ln\left(\frac{1+\alpha}{\alpha}\right)$$

e.g.

$$\alpha = 0.5 \text{ then } I_0 = 1.098612288668 \dots$$

$$\alpha = 2.0 \text{ then } I_0 = 0.405465108108 \dots$$

Step 2

Notice:

$$\begin{aligned}
 I_{n+1} &= \int_0^1 \frac{x^{n+1}}{x+\alpha} dx = \int_0^1 \frac{x^n(x+\alpha-\alpha)}{x+\alpha} dx = \int_0^1 x^n dx - \alpha \int_0^1 \frac{x^n}{x+\alpha} dx \\
 I_{n+1} &= \frac{1}{n+1} - \alpha I_n
 \end{aligned}$$

$$\begin{aligned}
I_0 \\
I_1 &= 1 - \alpha I_0 \\
I_2 &= \frac{1}{2} - \alpha I_1 \\
&\vdots \\
I_{100} &= \frac{1}{100} - \alpha I_0
\end{aligned}$$

If $\alpha = 0.5$ then $I_{100} = 0.00664$
 if $\alpha = 2.0$ then $I_{100} = 2.1 * 10^{22}$

$$\|I_n\| \leq \frac{1}{1 + \alpha}$$

Let's analyze what is happening

Math:

$$I_0^{ex}, I_{n+1}^{ex} = \frac{1}{n+1} - \alpha I_n^{ex}$$

CS:

$$I_0^{app}, I_{n+1}^{app} = \frac{1}{n+1} - \alpha I_n^{app}$$

At every step, there is some error:

$$\begin{aligned}
e_0 &= I_0^{ex} - I_0^{app} \\
e_n &= I_n^{ex} - I_n^{app}
\end{aligned}$$

Notice:

$$\begin{aligned}
e_{n+1} &= I_{n+1}^{ex} - I_{n+1}^{app} \\
&= \left(\frac{1}{n+1} - \alpha I_n^{ex} \right) - \left(\frac{1}{n+1} - \alpha I_n^{app} \right) \\
&= -\alpha I_n^{ex} + \alpha I_n^{app} \\
&= -\alpha (I_n^{ex} - I_n^{app}) \\
&= -\alpha e_n \\
&= (-\alpha)^{n+1} e_0
\end{aligned}$$

If $\|\alpha\| \leq 1$ then $\|e_n\| \rightarrow 0$ as $n \rightarrow \infty$

If $\|\alpha\| \geq 1$ then $\|e_n\| \rightarrow \infty$ as $n \rightarrow \infty$

What to do when $\|\alpha\| \geq 1$:

$$\begin{aligned}
I_n &= \frac{1}{\alpha(n+1)} - \frac{1}{\alpha} I_{n+1} \\
e_n &= -\frac{1}{\alpha} e_{n+1}
\end{aligned}$$

If $\|\alpha\| \geq 1$ then work backwards, e.g. I_{100} , Do $I_{200}, I_{199} \dots$

2 Interpolation

Given n points $(x_1, y_1), \dots, (x_N, y_N)$ x_i distinct $x_1 < x_2 < \dots < x_N$
 $y = p(x)$, p should be 'nice'
 'nice':

- polynomial; piecewise polynomial;

2.1 Polynomial Interpolation

Given n points (x_i, y_i) $i = 1, 2, \dots, n$
 Find a polynomial having degree $< n$ satisfying $p(x_i) = y_i$

Example

$(-1, 3), (1, 1), (2, 2)$

$$\begin{aligned} p(x) &= c_0 + c_1x + c_2x^2 \\ p(-1) &= c_0 - c_1 + c_2 = 3 \\ p(1) &= c_0 + c_1 + c_2 = 1 \\ p(2) &= c_0 + 2c_1 + 4c_2 = 2 \end{aligned}$$

$$\begin{bmatrix} 1 & -1 & 1 & 3 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 2 \end{bmatrix}$$

$$p(x) = \frac{4}{3} - x + \frac{2}{3}x^2$$

Given n points $(x_1, y_1), \dots, (x_N, y_N)$ x_i distinct $x_1 < x_2 < \dots < x_N$

- 1) Does there exist a polynomial $p(x)$ of degree $< n$ which interpolate the n points?
- 2) If it exists then is it unique?

$$\begin{aligned} &(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \\ p(x) &= c_1 + c_2x + \dots + c_nx^{n-1} \end{aligned}$$

n unknowns
 n equations

$$\begin{aligned} c_1 + c_2x_1 + \dots + c_nx_1^{n-1} &= y_1 \\ c_1 + c_2x_2 + \dots + c_nx_2^{n-1} &= y_2 \\ &\vdots \\ c_1 + c_2x_n + \dots + c_nx_n^{n-1} &= y_n \end{aligned}$$

Vandermonde matrix

$$V \cdot \vec{c} = \vec{y}$$

$$\det(V) = \prod_{i>j} (x_i - x_j) \neq 0 \quad \text{since all } x_i \text{ distinct}$$

2.2 Lagrange Form of Interpolating Polynomial

$$\begin{aligned} p(x) &= c_1 + c_2x + c_3x^2 + \dots + c_nx^{n-1} \\ &= y_1L_1(x) + y_2L_2(x) + \dots + y_nL_n(x) \end{aligned}$$

where each $L_i(x)$ is a polynomial of degree $< n$ which satisfies $L_i(x_i) = 1, L_i(x_j) = 0$ if $i \neq j$
 $L_i(x) \equiv$ Lagrange polynomial

$$L_i(x) = \frac{(x - x_1) \dots (x - x_{i-1})(x - x_{i+1}) \dots (x - x_n)}{(x_i - x_1) \dots (x_i - x_{i-1})(x_i - x_{i+1}) \dots (x_i - x_n)}$$

2.3 Hermite Interpolation

Points: $(x_L, y_L), (x_R, y_R)$

Derivative values: S_L, S_R

$S(x)$ degree $< 4, S(X_L) = Y_L, S(X_R) = Y_R, S'(X_L) = S_L, S'(X_R) = S_R$, Find $S(X)$

$$S(x) = c_1 + c_2(x - x_L) + c_3(x - x_L)^2 + c_4(x - x_L)^3$$

$$\begin{bmatrix} 1 & 0 & 0 & 0 & y_L \\ 0 & 1 & 0 & 0 & s_L \\ 1 & \Delta x & \Delta x^2 & \Delta x^3 & y_R \\ 0 & 1 & 2\Delta x & 3\Delta x^2 & s_R \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 0 & 0 & 2 \\ 0 & 1 & 0 & 0 & -1 \\ 1 & 2 & 4 & 8 & 4 \\ 0 & 1 & 4 & 12 & -1 \end{bmatrix}$$

$$y'_L = \frac{y_R - y_L}{x_R - x_L} = \text{slope}$$

$$c_1 = y_L$$

$$c_2 = s_L$$

$$c_3 = \frac{3y'_L - 2s_L - s_R}{\Delta x}$$

$$c_4 = \frac{s_L + s_R - 2y'_L}{\delta x^2}$$

Point Polynomial interpolation is good at the interpolating points but might not be very good elsewhere

We will need to try other methods to fit data, cubic splines

2.4 Cubic Splines

Given N points: (x_i, y_i)

A cubic spline is a function $S(x)$ which satisfies the following conditions

1. in each interval $[x_i, x_{i+1}]$, $S(x) = S_i(x)$ is a polynomial of degree at most 3.
2. $S(x)$ interpolates the points $(x_1, y_1), \dots, (x_N, y_N)$
3. $S'(x)$ exists and is continuous everywhere in $[x_1, x_N]$
4. $S''(x)$ exists and is continuous everywhere in $[x_1, x_N]$
5. 2 Boundary Conditions

As it stands as cubic spline problem has A unknowns and B equations, We want $A = B$

How many unknowns?

4 unknowns per interval, $N - 1$ intervals $\Rightarrow 4(N - 1) = 4N - 4$ unknowns

How many equations?

Condition (2): 2 per interval $\Rightarrow 2(N - 1) = 2N - 2$

Condition (3): 1 per interior pt $N - 2$

Condition (4): 1 per interior pt $N - 2$

$2N - 2 + N - 2 + N - 2 = 4N - 6$

There fore need 2 more conditions

Typical Boundary Conditions

1. Natural spline $S''(x_1) = 0, S''(x_N) = 0$
2. Clamped spline $S'(x_1) = s_1, S'(x_N) = s_n, s_1, s_n$ given
3. periodic spline $S'(x_1) = S'(x_N), S''(x_1) = S''(x_N)$
4. Not-a-knot (Matlab default) S''' continuous at x_2, x_{N-1}

How to compute a cubic spline?

1. Set a system of $4N - 4$ equations in the $4N - 4$ unknowns

$$S_i(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3$$

Solving Cost: $O(N^3)$ We want a method which compute a cubic spline using $O(N)$ operations

Set up a linear system for the unknown derivative values s_1, s_2, \dots, s_n which will be fast to solve and which will give us $S_1(x), \dots, S_{N-1}(x)$