

## FP Function Approximations for $\log_2()$ , $2^x$ , & $\text{powf}(x, r)$ Routines

First a review of how a single precision (32-bit) floating point number is represented in memory:

Bit Position																																		
31	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0			
$s_x$	$e_x = e[7:0] = \sum_{i=0}^7 e[i] \cdot 2^i$								$m_x = m[22:0] = \sum_{i=0}^{22} m[i] \cdot 2^i$																									
Integer Interpretation									$X = s_x \cdot 2^{31} + e_x \cdot 2^{23} + m_x$																									
Floating Point Interpretation									$x = (-1)^s \left(1 + \frac{m_x}{2^{23}}\right) \cdot 2^{e_x-127} = (-1)^s \cdot m'_x 2^{e'_x}$ $m'_x = \left(1 + \frac{m_x}{2^{23}}\right), \quad e'_x = (e_x - 127), \quad 1 \leq m'_x < 2, \quad -127 \leq e'_x \leq 128$																									
Hexadecimal Floating-Point Representation									Is a compact method of representing floating point values that prevents improper decimal truncation. Supported by C99 spec.  $\pi = 3.1415927410125732421875 = 0x3.243F6Cp0$ $e^1 = 2.71828174591064453125 = 0x2.B7E15p0$ $\frac{1}{\sqrt{2 \cdot \pi}} = 0.398942291736602783203125 = 0x6.62115p-4$																									

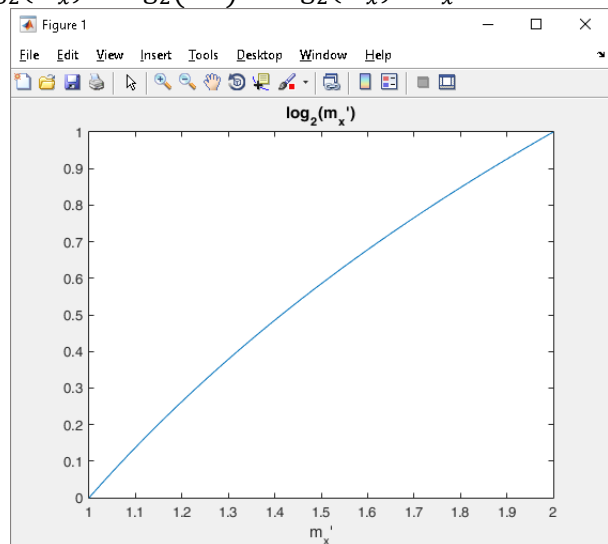
There are a couple caveats to representing particular numbers:

$x$	$s_x$	$e_x$	$m_x$
0	X	0	0
$+\infty$	0	255	0
$-\infty$	1	255	0

### Approximating the logarithm base 2 function

$$y = \log_2|x|, \quad y = \log_2(m'_x \cdot 2^{e'_x}) = \log_2(m'_x) + \log_2(2^{e'_x}) = \log_2(m'_x) + e'_x$$

Calculating the function approximation via range reduction allows for much higher precision polynomial fit over such a short interval. (i.e.  $1 \leq m'_x < 2$ )



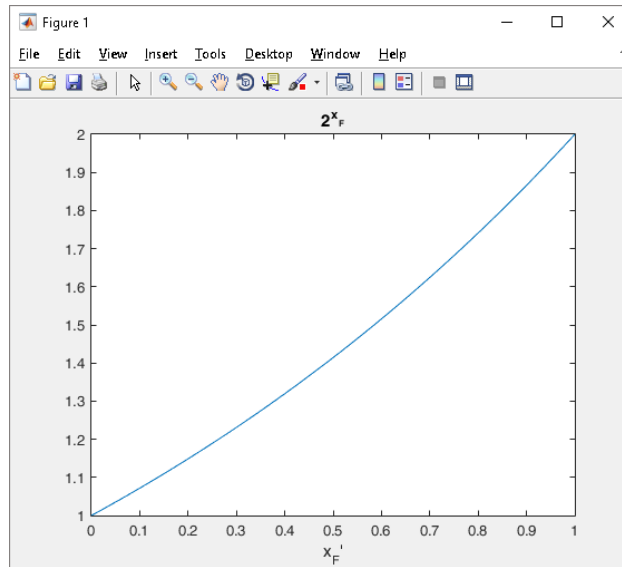
Using an open source Linux package Sollya to calculate a bounded N-th order polynomial

## FP Function Approximations for $\log_2()$ , $2^x$ , & $\text{powf}(x, r)$ Routines

### Approximating the exponent base 2 function

$$y = 2^x, \quad m'_y \cdot 2^{e'_y} = 2^x = 2^{x_I + x_F}, \quad \text{where } x_I = \lfloor x \rfloor, \quad x_F = x - \lfloor x \rfloor \text{ then}$$

$$m'_y \cdot 2^{e'_y} = 2^{x_I + x_F} = 2^{x_F} \cdot 2^{x_I}, \quad e'_y = x_I, \quad m'_y = 2^{x_F}, \quad 0 \leq x_F < 1, \quad 1 \leq m'_y < 2$$



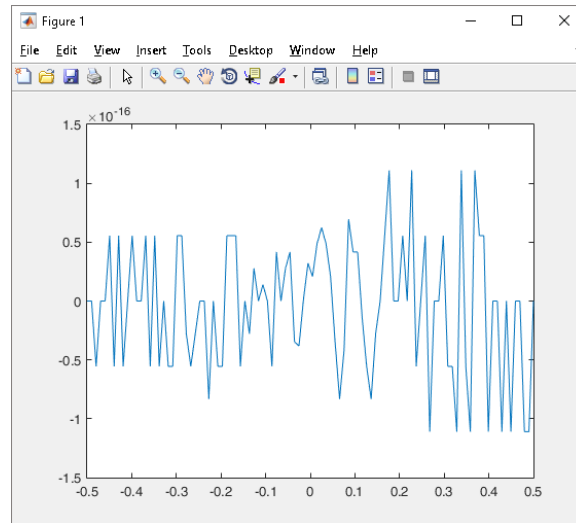
It's important to note that the approximation nodes must be exact at 0 & 1 which requires a hand tuned polynomial approximation using Newton's Divided Difference or Lagrange Interpolation method to force both end-points to be exact. **Caution** this method has an inherent flaw, the final error from the bounded polynomial grows with the size of the input  $x$  exponentially. The only way to counter this would be to use a pade approximation such as the one used in the boost library:

$$y = x \left( c_0 + \frac{c_{n0} + c_{n1} \cdot x^1 + c_{n2} \cdot x^2 + c_{n3} \cdot x^3 + c_{n4} \cdot x^4 + c_{n5} \cdot x^5}{c_{d0} + c_{d1} \cdot x^1 + c_{d2} \cdot x^2 + c_{d3} \cdot x^3 + c_{d4} \cdot x^4 + c_{d5} \cdot x^5} \right), \quad -\frac{1}{2} < x < \frac{1}{2}$$

$$c_0 = 0.10281276702880859e1$$

$$c_n = \begin{cases} -0.28127670288085937e-1 & 0.51278186299064534e0 & -0.6310029069350198e-1 & 0.11638457975729296e-1 \end{cases}$$

$$c_d = \begin{cases} -0.52143390687521003e-3 & 0.21491399776965688e-4 \\ 1.0 & -0.45442309511354755e0 & 0.90850389570911714e-1 & -0.10088963629815502e-1 \\ 0.63003407478692265e-3 & -0.17976570003654402e-4 \end{cases}$$



This pade approximation introduces increased complexity as well as division which is bad news bears for decreasing functional latency. It is also not necessary as we will see the accuracy is adequate if we restrict the input  $0 < x \leq 1$

## FP Function Approximations for $\log_2()$ , $2^x$ , & $\text{powf}(x, r)$ Routines

### Approximating the power function

$$y = x^p, \quad \log_2(y) = \log_2(x^p) = p \cdot \log_2(x), \quad y = 2^{p \cdot \log_2(x)}, \quad \text{for } 0 < x \leq 1, \quad -\infty < p \cdot \log_2(x) \leq 0$$

So the exponent will always be negative. A special restriction must be applied to prevent  $\log_2(0)$  and  $-127 \leq p \cdot \log_2(x)$ . Below is functional error for all combinations of possible ordered polynomials for both the  $\text{exp2}(x)$  and  $\log_2(x)$ .

