

A vertical grey bar is on the left side of the slide. An orange arrow-shaped box points to the right from the bar, containing the date. Below the bar, several thin, curved lines in black and grey extend upwards and to the right.

November 17, 2013

Computation of Exponentials and Logarithms

With binary arithmetics

Jarmo Nikkanen

jnikkanen@kyp.net

THE LICENSE

Copyright © 2013 Jarmo Nikkanen

Permission is hereby granted, free of charge, to any person obtaining a copy of this software and associated documentation (the "Software"), to deal in the Software without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Software.

THE SOFTWARE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE SOFTWARE OR THE USE OR OTHER DEALINGS IN THE SOFTWARE.

FOREWORD

This document is created in a hope that it might give some ideas on how to compute 2-based exponentials and logarithms. Source codes are provided because some people can read code better than a mathematical descriptions. Remember that the codes are provided to clarify the idea, NOT to provide a well working solution. Source codes are missing rounding of the result, special case handling and additional shifts to improve precision.

2-BASED EXPONENTIAL FUNCTION WITHOUT MULTIPLICATION

Index	Multiplier	fSub	dwSub
n	$m = 1 + 1/2^n$	$s = \log_2 m$	
0	2.0000000	1.0000000	0x800000
1	1.5000000	0.5849625	0x4AE00D
2	1.2500000	0.3219281	0x2934F0
3	1.1250000	0.1699250	0x15C01A
4	1.0625000	0.0874628	0x0B31FB
5	1.0312500	0.0443941	0x05AEB4
6	1.0156250	0.0223678	0x02DCF2
7	1.0078125	0.0112273	0x016FE5
...

The method described in a previous section could be good one if you have a fast hardware multiplier available, otherwise the multiplication might become a problem. There is a way to avoid the multiplication. As you probably already know that a multiplication or a division with a powers of two can be completed by shifting the bits. So, we could try to approach the problem from the opposite direction by trying to find the multipliers those will allow us to do the multiplication easily. One multiplier that fits in our needs is $[1 + 1/2^n]$. If we want to multiply a value x with this multiplier then we can do it by shifting the bits: $x = x + (x \gg n)$; But in that case we need a table of corresponding subtraction values. We can't use this multiplier in a per bit basis as we did in the previous

section. The subtraction values can be computed by taking a base-2 logarithm from the multiplier. The first eight values are shown in a table right here. Upper 8-bits of dwSub are unused in this example, it would be a good idea to put them in a good use.

```
float Exp2(float iv)
{
    int i = 1;
    int e = Exponent(iv);
    DWORD m = Mantissa(iv) | 0x800000;
    DWORD b = Shift(m, e) & 0x7FFFFFFF; // Shift bits to left or right based on polarity
    DWORD v = 0x800000; // Initial starting value

    while (i<23) {
        if (b>=dwSub[i]) {
            b-=dwSub[i];
            v = v + (v>>i);
        } else i++;
    }

    v &= 0x7FFFFFFF; // Form a floating point number. Remove implicit bit.
    v += 127<<23; // Initialize exponent
    if (e>=0) v += (m>>(23-e))<<23; // Compute the exponent and add it to the float
    return *(float*)&v; // Convert dword to float
}
```

The code above is simplified to use only a positive values. Also, you may need to adjust the fixed point by bit shifting to reduce a round-off error accumulation. These shifts are removed from the code as well as rounding of the result. To compute the exponential using a negative input value you need to use $[1 - 1/2^n]$ multiplier instead of $[1 + 1/2^n]$. See the 2-based logarithm section for some additional details.

2-BASED LOGARITHMIC FUNCTION

Index	Bit weight	fDiv	iDiv
b	$w = 1/2^{23-b}$	$m = 2^w$	$i = 1/2^w$
22	0.5000000	1.4142135	0.7071068
21	0.2500000	1.1892071	0.8408964
20	0.1250000	1.0905077	0.9170040
19	0.0625000	1.0442737	0.9576033
18	0.0312500	1.0218972	0.9785720
17	0.0156250	1.0108893	0.9892280
16	0.0078125	1.0054299	0.9945995
15	0.0039063	1.0027113	0.9972960
...

The computation process of 2-based logarithm is pretty much the opposite than with the exponential. As long as the source (input) value is greater than 1.0. The input value is divided by a known constant (as large as possible but less than the source value) and then a 2-based logarithm of the constant is added to the output value. Each division of the source value will make it to approach 1.0. In this example case adding of $\log_2 c[i]$ is simplified to a bit setting. Here is a table containing a divisor and it's inverse. The values are the same as in the first section. We could have a longer table containing a divisors for bits 23-31 but it would be a waste. Instead, we can simply divide the input value with 2^E in other words setting the exponent E to 127 from

the source value (that is zero due to exponent bias). That will scale the input value in range [1-2]. Then we have to add the original value of the exponent into the final result. The exponent E of the input value itself is the integer (whole number) part of the final result. In the code below we subtract 1 from the 'e', because the implicit bit is added to the float returned by CreateFloat().

```
float Log2(float iv)
{
    DWORD a = 0;
    int e = Exponent(iv);
    SetExponent(&iv, 0); // Set exponent to zero (i.e. scale input range to [1-2])
    for (int i=22; i>=1; i--) {
        if (iv>=fDiv[i]) iv*=iDiv[i], a|=1; // Divide and set lowest bit
        a<<=1; // Shift the bits
    }
    return Int2Float(e-1) + CreateFloat(a);
}
```

2-BASED LOGARITHM WITHOUT MULTIPLICATION

This section will describe another method using a bit shifting instead of floating point multiplication. The principles are exactly the

Index	Multiplier	fAdd	dwAdd
n	$m = 1 - 1/2^n$	$ \log_2 m $	
1	0.5000000	1.0000000	0x800000
2	0.7500000	0.4150375	0x351FF3
3	0.8750000	0.1926451	0x18A898
4	0.9375000	0.0931094	0x0BEB02
5	0.9687500	0.0458037	0x05DCE5
6	0.9843750	0.0227201	0x02E87D
7	0.9921875	0.0113153	0x0172C7
8	0.9960938	0.0056466	0x00B906
...

same as described in an earlier sections. Also the table that is used here is exactly the same as the one required by negative values earlier. It doesn't really matter what value is used in a division/multiplication as long as the source value approaches 1.0 in every step and the relationship between the divisor and the value being added is correct. Also, note that you may need to use the same divisor/multiplier more than once. For an example: if the input value is 1.9 you need to multiply it twice with 0.75. That is why the i++ is in the 'else' statement.

```
float Log2(float iv)
{
    DWORD lg = 0;
    int i = 2;
    int m = Mantissa(iv) | 0x800000;
    while (i<23) {
        int t = m - (m>>i);
        if (t>0x800000) {
            m = t;
            lg += dwAdd[i]; // Table index 0 is unused
        } else i++;
    }
    return Int2Float(Exponent(iv)-1) + CreateFloat(lg);
}
```

2-BASED LOGARITHM WITHOUT TABLE

If we look at the function we developed in the section [2-based logarithmic function](#). We can notice that the actual computation is pretty simple since the $fDiv[i] = \sqrt[x]{2}$, where $x = 2^i$.

$$\text{if } (v \geq \sqrt[x]{2}) \text{ then } v = v / \sqrt[x]{2}$$

We are comparing v to the x :th root of 2. Then, why not raise both sides of the “statement” to power of x , resulting following statement below. This can be easily archived by squaring the input value v each step in a loop. Also remember that division by two can be completed by decrementing the exponent E .

$$\text{if } (v^x \geq 2) \text{ then } v^x = v^x / 2$$

```
float Log2(float iv)
{
    DWORD a = 0;
    int e = Exponent(iv);
    SetExponent(&iv,0); // Set exponent to zero (i.e. scale input range to [1-2])
    for (int i=22;i>=1;i--) {
        iv*=iv;
        if (iv>=2.0f) iv*=0.5f, a|=1; // Set lowest bit
        a<<=1; // Shift the bits
    }
    return Int2Float(e-1) + CreateFloat(a);
}
```

2-BASED EXPONENTIAL WITHOUT TABLE

Computation of 2-based exponential without a table is pretty easy. We can simply replace the table by taking a square root from a previous root in a loop. But doing so would be mathematically expensive. Therefore, it’s not much an option.

```
float sq = 2.0f;
for (int i=0;i<23;i++) {
    sq=sqrt(sq); // Take a square root
    if (b&0x400000) rv *= sq; // Perform a multiplication for each active bit in a fixed point fraction
    b<<=1; // Shift the bits
}
```

In theory, it would be possible to take the n :th root of the lowest bit and start squaring it up, but the value is so close to 1.0 that a computer doesn’t really make a difference between the two. There isn’t enough precision to do the squaring without doubling the bits in the mantissa. The code below will show the idea.

```
float sq = 1.0000000826295864; // 2^23:th root of 2
for (int i=0;i<23;i++) {
    if (b&1) rv *= sq; // Perform a multiplication for each active bit in a fixed point fraction
    b>>=1; // Shift the bits
    sq*=sq; // Square it up
}
```

How about splitting the value 1.0000000826295864 into a two different values 1.0 and 8.26295864e-8. You probably recall from the school that $(a + b)^2 = a^2 + 2ab + b^2$. This will work in our favor. In our case the variable a would be 1.0 and that will simplify the equation to $(1 + b)^2 = 1 + 2b + b^2$. Also we don’t want to compute the square, instead we want the next b_{n+1} which is given by following equation below. *Last two underlined forms are usable.*

$$b_{n+1} = (1 + b_n)^2 - 1 = (1 + 2b_n + b_n^2) - 1 = \underline{2b_n + b_n^2} = \underline{b_n(2 + b_n)}$$

```
float Exp2(float iv)
{
    DWORD m = Mantissa(iv) | 0x800000;
    int e = Exponent(iv);
    DWORD b = Shift(m, e); // Shift bits to left or right based on polarity
    float rv = 1.0f;
    float sq = 8.26295864e-8f; // (2^23:th root of 2) minus 1.0

    for (int i=0;i<23;i++) {
        if (b&1) rv *= (1.0f + sq); // Perform a multiplication for each active bit
        sq *= (2.0f+sq); // Square it up
        b>>=1; // Shift the bits
    }
    DWORD *d = (DWORD*)&rv;
    if (e>=0) *d += (m>>(23-e))<<23; // Add the integer part to the exponent
    return rv;
}
```

SOME BASIC ALGORITHMS

```
int Exponent(float f)
{
    DWORD *d = (DWORD*)&f;
    int e = (*d>>23)&0xFF;
    return e - 127; // Return signed unbiased exponent
}

DWORD Mantissa(float f)
{
    DWORD *d = (DWORD*)&f;
    return *d&0x7FFFFFFF;
}

bool Sign(float f)
{
    DWORD *d = (DWORD*)&f;
    return ((*d&0x80000000)!=0);
}

float CreateFloat(DWORD m, int e, bool s)
{
    DWORD d = (m&0x7FFFFFFF) + ((e+127)<<23);
    if (s) d |= 0x80000000; // Add a sign bit
    return *(float*)&d;
}

float CreateFloat(DWORD m)
{
    DWORD d = (m&0x7FFFFFFF) + (127<<23);
    return *(float*)&d;
}

void SetExponent(float *f, int e)
{
    *f = CreateFloat(Mantissa(*f), e, Sign(*f));
}

DWORD Shift(DWORD w, int d)
{
    if (d<0) return w>>(-d);
    return w<<d;
}

DWORD HighestBitZero(DWORD *x)
{
    for (int i=0;i<32;i++) {
        DWORD q = 0x80000000>>i;
        if (*x&q) { *x-=q; return 31-i; }
    }
    return 0;
}
```