

C++ Implementation

Pieter P

Dividing by Powers of 2

The factor α in the difference equation of the Exponential Moving Average filter is a number between zero and one. There are two main ways to implement this multiplication by α : Either we use floating point numbers and calculate the multiplication directly, or we use integers, and express the multiplication as a division by $1/\alpha > 1$.

Both floating point multiplication and integer division are relatively expensive operations, especially on embedded devices or microcontrollers.

We can, however, choose the value for α in such a way that $1/\alpha = 2^k, k \in \mathbb{N}$.

This is useful, because a division by a power of two can be replaced by a very fast right bitshift:

$$\alpha \cdot x = \frac{x}{2^k} = x \gg k$$

We can now rewrite the difference equation of the EMA with this optimization in mind:

$$\begin{aligned} y[n] &= \alpha x[n] + (1 - \alpha)y[n - 1] \\ &= y[n - 1] + \alpha(x[n] - y[n - 1]) \\ &= y[n - 1] + \frac{x[n] - y[n - 1]}{2^k} \\ &= y[n - 1] + (x[n] - y[n - 1]) \gg k \end{aligned}$$

Negative Numbers

There's one caveat though: this doesn't work for negative numbers. For example, if we try to calculate the integer division $-15/4$ using this method, we get the following answer:

$$\begin{aligned} -15/4 &= -15 \cdot 2^{-2} \\ -15 \gg 2 &= 0b11110001 \gg 2 \\ &= 0b11111100 \\ &= -4 \end{aligned}$$

This is not what we expected! Integer division in programming languages such as C++ returns the quotient truncated towards zero, so we would expect a value of -3 . The result is close, but incorrect nonetheless.

This means we'll have to be careful not to use this trick on any negative numbers. In our difference equation, both the input $x[n]$ and the output $y[n]$ will generally be positive numbers, so no problem there, but their difference can be negative. This is a problem. We'll have to come up with a different representation of the difference equation that doesn't require us to divide any negative numbers:

$$\begin{aligned} y[n] &= y[n - 1] + \alpha(x[n] - y[n - 1]) \\ y[n] &= y[n - 1] + \frac{x[n] - y[n - 1]}{2^k} \\ 2^k y[n] &= 2^k y[n - 1] + x[n] - y[n - 1] \\ z[n] &\triangleq 2^k y[n] \Leftrightarrow y[n] = 2^{-k} z[n] \\ z[n] &= z[n - 1] + x[n] - 2^{-k} z[n - 1] \end{aligned}$$

We now have to prove that $z[n - 1]$ is greater than or equal to zero. We'll prove this using induction:

Base case: $n - 1 = -1$

The value of $z[-1]$ is the initial state of the system. We can just choose any value, so we'll pick a value that's greater than or equal to zero: $z[-1] \geq 0$.

Induction step: n

Given that $z[n - 1] \geq 0$, we can now use the difference equation to prove that $z[n]$ is also greater than zero:

$$z[n] = z[n - 1] + x[n] - 2^{-k} z[n - 1]$$

We know that the input $x[n]$ is always zero or positive.

Since $k > 1 \Rightarrow 2^{-k} < 1$, and since $z[n - 1]$ is zero or positive as well, we know that

$$z[n - 1] \geq 2^{-k} z[n - 1] \Rightarrow z[n - 1] - 2^{-k} z[n - 1] \geq 0.$$

Therefore, the entire right-hand side is always positive or zero, because it is a sum of two numbers that are themselves greater than or equal to zero. \square

Rounding

A final improvement we can make to our division algorithm is to round the result to the nearest integer, instead of truncating it towards zero.

Consider the rounded result of the division a/b . We can then express it as a flooring of the result plus one half:

$$\begin{aligned}\left\lceil \frac{a}{b} \right\rceil &= \left\lceil \frac{a}{b} + \frac{1}{2} \right\rceil \\ &= \left\lceil \frac{a + \frac{b}{2}}{b} \right\rceil\end{aligned}$$

When b is a power of two, this is equivalent to:

$$\begin{aligned}\left\lceil \frac{a}{2^k} \right\rceil &= \left\lceil \frac{a}{2^k} + \frac{1}{2} \right\rceil \\ &= \left\lceil \frac{a + \frac{2^k}{2}}{2^k} \right\rceil \\ &= \left\lceil \frac{a + 2^{k-1}}{2^k} \right\rceil \\ &= (a + 1 \ll (k - 1)) \gg k\end{aligned}$$

Implementation in C++

We now have everything in place to write an implementation of the EMA in C++:

EMA.cpp

```
1  #include <stdint.h>
2
3  template <uint8_t K, class uint_t = uint16_t>
4  class EMA {
5  public:
6      uint_t operator()(uint_t x) {
7          z += x;
8          uint_t y = (z + (1 << (K - 1))) >> K;
9          z -= y;
10         return y;
11     }
12
13     static_assert(
14         uint_t(0) < uint_t(-1), // Check that `uint_t` is an unsigned type
15         "Error: the uint_t type should be an unsigned integer, otherwise, "
16         "the division using bit shifts is invalid.");
17
18     private:
19         uint_t z = 0;
20     };
```

Note how we save $z[n] - 2^{-k}z[n]$ instead of just $z[n]$. Otherwise, we would have to calculate $2^{-k}z[n]$ twice (once to calculate $y[n]$, and once on the next iteration to calculate $2^{-k}z[n - 1]$), and that would be unnecessary.

Signed Rounding Division

It's possible to implement a signed division using bit shifts as well. The only difference is that we have to subtract 1 from the dividend if it's negative.

On ARM and x86 platforms, the performance difference between the signed and unsigned version is small.

On some other architectures, like the AVR architecture used by some Arduino microcontrollers, the signed version is significantly slower.

I provided two implementations of the signed division. Notice how on x86 and ARM the second one is faster, while on AVR, the first one is faster.

The code was compiled using the `-O2` optimization level.

Implementation of Signed and Unsigned Division by a Multiple of Two

```

1  constexpr unsigned int K = 3;
2
3  signed int div_s1(signed int val) {
4      int round = val + (1 << (K - 1));
5      if (val < 0)
6          round -= 1;
7      return round >> K;
8  }
9
10 signed int div_s2(signed int val) {
11     int neg = val < 0 ? 1 : 0;
12     return (val + (1 << (K - 1)) - neg) >> K;
13 }
14
15 unsigned int div_u(unsigned int val) {
16     return (val + (1 << (K - 1))) >> K;
17 }

```

Assembly Generated on x86_64 (GCC 9.2)

```

1  div_s1(int):
2      mov     eax, edi
3      not     eax
4      shr     eax, 31
5      lea     eax, [rax+3+rdi]
6      sar     eax, 3
7      ret
8
9  div_s2(int):
10     lea     eax, [rdi+4]
11     shr     edi, 31
12     sub     eax, edi
13     sar     eax, 3
14     ret
15
16 div_u(unsigned int):
17     lea     eax, [rdi+4]
18     shr     eax, 3
19     ret

```

Assembly Generated on ARM 64 (GCC 8.2)

```

1  div_s1(int):
2      mvn     w1, w0
3      add     w0, w0, w1, lsr 31
4      add     w0, w0, 3
5      asr     w0, w0, 3
6      ret
7
8  div_s2(int):
9      add     w1, w0, 4
10     sub     w0, w1, w0, lsr 31
11     asr     w0, w0, 3
12     ret
13
14 div_u(unsigned int):
15     add     w0, w0, 4
16     lsr     w0, w0, 3
17     ret

```

Assembly Generated on AVR (GCC 5.3)

```

1  __zero_reg__ 1
2  div_s1(int):
3      sbrc r25,7 # Skip if Bit in Register is Cleared: val < 0
4      rjmp .L2
5      # val >= 0
6      adiw r24,4 # Add Immediate to Word: val + (1 << (K - 1)) = val + 4
7      asr r25 # Arithmetic Shift Right: shift high byte (preserve sign)
8      ror r24 # Rotate Right through Carry: shift low byte
9      asr r25 # Two more times
10     ror r24
11     asr r25
12     ror r24
13     ret
14 .L2:
15     # val < 0
16     adiw r24,3 # Add Immediate to Word: val + (1 << (K - 1)) - 1 = val + 3
17     asr r25 # Arithmetic Shift Right: shift high byte (preserve sign)
18     ror r24 # Rotate Right through Carry: shift low byte
19     asr r25 # Two more times
20     ror r24
21     asr r25
22     ror r24
23     ret
24
25 div_s2(int):
26     movw r18,r24
27     subi r18,-4 # Subtract immediate: val + (1 << (K - 1)) = val + 4
28     sbci r19,-1 # Subtract Immediate with Carry: (low byte)
29     mov r24,r25
30     rol r24 # Rotate Left through Carry: C flag is now sign bit
31     clr r24 # Clear Register: set 24 to 0
32     rol r24 # Rotate Left through Carry: lsb is now sign bit
33     movw r20,r18
34     sub r20,r24 # Subtract without Carry: val + 4 - neg
35     sbc r21,__zero_reg__ # Subtract with Carry: (low byte)
36     movw r24,r20
37     asr r25 # Arithmetic Shift Right: shift high byte (preserve sign)
38     ror r24 # Rotate Right through Carry: shift low byte
39     asr r25 # Two more times
40     ror r24
41     asr r25
42     ror r24
43     ret
44
45 div_u(unsigned int):
46     adiw r24,4 # Add Immediate to Word: val + (1 << (K - 1)) = val + 4
47     lsr r25 # Logical Shift Right: shift high byte (no sign extension)
48     ror r24 # Rotate Right through Carry: shift low byte
49     lsr r25 # Two more times
50     ror r24
51     lsr r25
52     ror r24
53     ret

```

Keep in mind that an `int` on AVR is only 16 bits wide, whereas an `int` on ARM or x86 is 32 bits wide. If you use 32-bit integers on AVR, the result is even more atrocious.

You can try it for yourself on the [Compiler Explorer](#).

Arduino Example

```

1  template <uint8_t K, class uint_t = uint16_t>
2  class EMA {
3  public:
4      uint_t operator()(uint_t x) {
5          z += x;
6          uint_t y = (z + (1 << (K - 1))) >> K;
7          z -= y;
8          return y;
9      }
10
11      static_assert(
12          uint_t(0) < uint_t(-1), // Check that `uint_t` is an unsigned type
13          "Error: the uint_t type should be an unsigned integer, otherwise, "
14          "the division using bit shifts is invalid.");
15
16  private:
17      uint_t z = 0;
18  };
19
20  void setup() {
21      Serial.begin(115200);
22      while (!Serial);
23  }
24
25  const unsigned long interval = 10000; // 100 Hz
26
27  void loop() {
28      static EMA<2> filter;
29      static unsigned long prevMicros = micros();
30      if (micros() - prevMicros >= interval) {
31          int rawValue = analogRead(A0);
32          int filteredValue = filter(rawValue);
33          Serial.print(rawValue);
34          Serial.print('\t');
35          Serial.println(filteredValue);
36          prevMicros += interval;
37      }
38  }

```