

Simple Moving Average

Pieter P

Difference equation

The difference equation of the Simple Moving Average filter is derived from the mathematical definition of the average of N values: the sum of the values divided by the number of values.

$$y[n] = \frac{1}{N} \sum_{i=0}^{N-1} x[n-i]$$

In this equation, $y[n]$ is the current output, $x[n]$ is the current input, $x[n-1]$ is the previous input, etc. N is the length of the average.

Impulse and step response

From the previous equation, we can now easily calculate the impulse and step response.

The impulse response is the output of the filter when a Kronecker delta function is applied to the input. Recall the definition of the Kronecker delta:

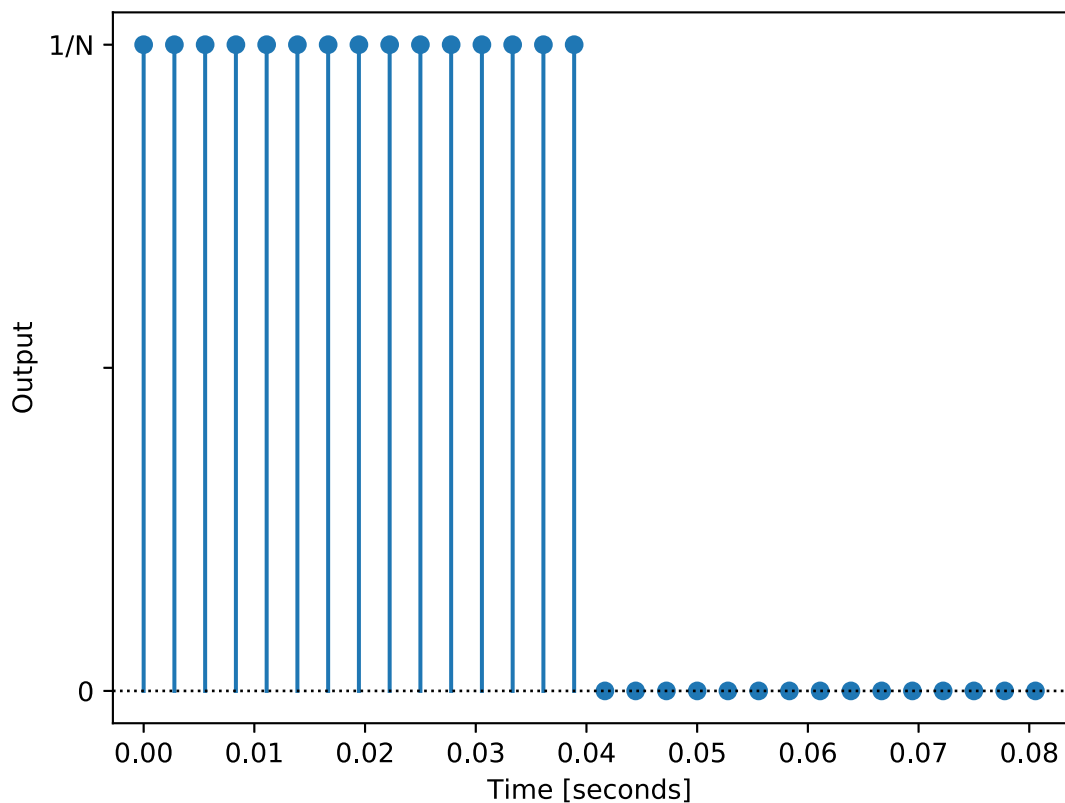
$$\delta[n] = \begin{cases} 1 & n = 0 \\ 0 & n \neq 0. \end{cases}$$

The impulse response of the SMA is

$$\begin{aligned} y_{\text{impulse}}[n] &= h[n] \\ &= \frac{1}{N} \sum_{i=0}^{N-1} \delta[n-i] \\ &= \begin{cases} 1/N & 0 \leq n < N \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

For example, if $N = 15$, the impulse response is as follows:

Impulse Response



The step response is the output of the filter when a Heaviside step function is applied to the input. The Heaviside step function is defined as

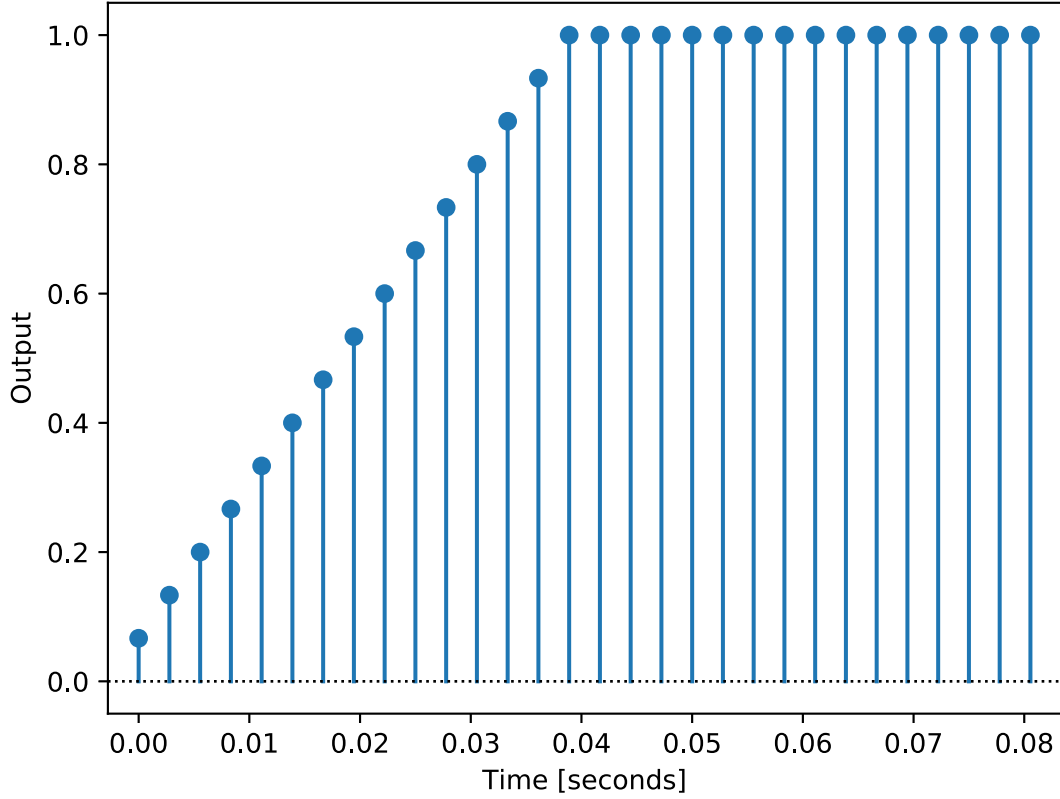
$$H[n] = \begin{cases} 0 & n < 0 \\ 1 & n \geq 0. \end{cases}$$

The step response of the SMA is

$$y_{step}[n] = \begin{cases} 0 & n < 0 \\ (n+1)/N & 0 \leq n < N \\ 1 & n \geq N. \end{cases}$$

For example, if $N = 15$, the step response is as follows:

Step Response



Transfer function

The output of discrete-time linear time-invariant (DTLTI) systems, of which the SMA is an example, can be expressed as the convolution of the input with the impulse response. In other words, the impulse response describes the behavior of the system, for all possible inputs.

To prove the previous statement, we'll start with the following trivial property: any signal $x[n]$ can be expressed as a convolution of the Kronecker delta function with itself, that is

$$x[n] = x[n] * \delta[n] = \sum_{k=-\infty}^{+\infty} x[k] \delta[n - k].$$

You can easily see that all terms where $k \neq n$ are zero, because the Kronecker delta is zero in that case. Only the term for $k = n$ is non-zero, in which case the Kronecker delta is one, so the result is just $x[n]$.

You can also interpret this as the signal being made up of a sum of infinitely many scaled and shifted Kronecker delta functions.

Let T be the transformation performed by the DTLTI system, then $y[n]$ is the output after applying T to the input signal $x[n]$. The following derivation makes use of the fact that T is a linear transformation and that it is time-invariant:

$$\begin{aligned}
y[n] &= T(x[n]) \\
&= T\left(\sum_{k=0}^{\infty} x[k] \delta[k-n]\right) \\
&= \sum_{k=0}^{\infty} T(x[k] \delta[k-n]) \\
&= \sum_{k=0}^{\infty} x[k] T(\delta[k-n]) \\
&= \sum_{k=0}^{\infty} x[k] h[k-n] \\
&= x[n] * h[n]
\end{aligned}$$

Since the factor $x[k]$ is independent of time n , it can be moved outside of the T operator. T applied to the Kronecker delta is (by definition) the impulse response of T , $h[n]$, but shifted by k time steps.

Analysis of such systems is usually easier in the Z-domain, in which the convolution is reduced to a simple product. The (unilateral) Z-transform is defined as:

$$\mathcal{Z}\{x[n]\} = \sum_{n=0}^{\infty} x[n]z^{-n}$$

If $X(z) = \mathcal{Z}\{x[n]\}$, $Y(z) = \mathcal{Z}\{y[n]\}$ and $H(z) = \mathcal{Z}\{h[n]\}$ are the Z-transforms of the input, output and impulse response respectively, then:

$$\begin{aligned}
\mathcal{Z}\{y[n]\} &= \mathcal{Z}\{x[n] * h[n]\} \\
Y(z) &= \sum_{n=0}^{\infty} x[n] * h[n] z^{-n} \\
&= \sum_{n=0}^{\infty} \sum_{k=0}^{\infty} x[k] h[n-k] z^{-n} \\
&= \sum_{k=0}^{\infty} \sum_{n=0}^{\infty} x[k] h[n-k] z^{-n+k-k} \\
&= \sum_{k=0}^{\infty} x[k] z^{-k} \sum_{n=0}^{\infty} h[n-k] z^{-(n-k)} \\
Y(z) &= X(z)H(z) \\
H(z) &= \frac{Y(z)}{X(z)}
\end{aligned}$$

$H(z)$ is called the transfer function of the system.

Let's calculate the transfer function of the SMA.

We can use one of two approaches: use the difference equation and use some of the properties of the Z-transform to calculate $\frac{Y(z)}{X(z)}$, or apply the definition of the Z-transform to the impulse response $h[n]$ directly.

Using the difference equation

All you have to do is apply the time shifting property of the Z transform:

$$\forall n_0 \in \mathbb{N} : \mathcal{Z}\{y[n - n_0]\} = z^{-n_0} \mathcal{Z}\{y[n]\}$$

Then just use the fact that the Z-transform is linear, and the derivation is trivial:

$$\begin{aligned}
y[n] &= \frac{1}{N} \sum_{i=0}^{N-1} x[n-i] \\
\mathcal{Z}\{y[n]\} &= \mathcal{Z}\left\{\frac{1}{N} \sum_{i=0}^{N-1} x[n-i]\right\} \\
&= \frac{1}{N} \sum_{i=0}^{N-1} \mathcal{Z}\{x[n-i]\} \\
&= \frac{1}{N} \sum_{i=0}^{N-1} z^{-i} \mathcal{Z}\{x[n]\} \\
Y(z) &= \frac{1}{N} \sum_{i=0}^{N-1} z^{-i} X(z) \\
H(z) &\triangleq \frac{Y(z)}{X(z)} \\
&= \frac{1}{N} \sum_{i=0}^{N-1} z^{-i} \\
&= \frac{1}{N} \sum_{i=0}^{N-1} z^{-i+(N-1)-(N-1)} \\
&= \frac{1}{N} z^{-(N-1)} \sum_{i=0}^{N-1} z^{(N-1)-i} \\
&\quad j \triangleq N-1-i \\
&= \frac{1}{N} \frac{\sum_{j=0}^{N-1} z^j}{z^{N-1}} \\
&= \frac{1}{N} \frac{z^{N-1} + z^{N-2} + \dots + z^2 + z + 1}{z^{N-1}}
\end{aligned}$$

In the last steps, we just rearrange some terms and keep only positive powers of z in both the numerator and the denominator.

Using the impulse response

For this derivation, we can use the fact that the impulse response $h[n]$ is only non-zero for the first N terms.

$$\begin{aligned}
H(z) &= \mathcal{Z}\{h[n]\} \\
&= \sum_{n=0}^{\infty} h[n] z^{-n} \\
&= \sum_{n=0}^{N-1} \frac{1}{N} z^{-n}
\end{aligned}$$

It is clear that this expression is identical to the previous one.

We can solve the summation in a different way using the formula for the sum of a geometric series:

$$\sum_{n=0}^{N-1} r^n = \frac{1 - r^N}{1 - r}$$

The transfer function then becomes:

$$\begin{aligned}
H(z) &= \frac{1}{N} \sum_{n=0}^{N-1} z^{-n} \\
&= \frac{1}{N} \sum_{n=0}^{N-1} \left(\frac{1}{z}\right)^n \\
&= \frac{1}{N} \frac{1 - z^{-N}}{1 - z^{-1}} \\
&= \frac{1}{N} \frac{z^N - 1}{z^{N-1}(z - 1)}
\end{aligned}$$

In these expressions, z is a complex variable, and $H(z)$ is a complex function.

There are a couple of interesting values for z : values that result in the transfer function becoming zero, called *zeros*, and values that result in the transfer function becoming undefined, called *poles*.

They can be found as the roots of the numerator and the denominator respectively, after eliminating any common factors.

The roots of the numerator are the solutions of $z^N - 1 = 0$ or $z^N = 1$. Note that this equation has N solutions, namely all N -th roots of unity. One of the solutions is $z = 1$, but there are $N - 1$ complex solutions as well, evenly spaced along the unit circle: $e^{2\pi jn/N} \mid n \in [0, N - 1]$.

The roots of the denominator are the solutions of $z^{N-1}(z - 1) = 0$. This has just two distinct solutions: $z = 0$ with a multiplicity of $N - 1$ and $z = 1$ with a multiplicity of one.

Since both the numerator and the denominator have the root $z = 1$ in common, we have to factor out the common factor $(z - 1)$. This leaves us with the following $N - 1$ poles and zeros:

$$\begin{aligned} \text{zeros} &= \left\{ e^{2\pi jn/N} \mid n \in [1, N - 1] \right\} \\ \text{poles} &= \{0\}. \end{aligned}$$

Frequency response

An important property of discrete-time linear time-invariant systems is that it preserves the pulsantance (angular frequency) of sinusoidal signals, only the phase shift and the amplitude are altered. In other words, sines and cosines are eigenfunctions of DTLTI systems. This makes it relatively easy to express the frequency response (sometimes called magnitude response) of a filter.

We're interested in the spectrum of frequencies that a signal contains, so it makes sense to decompose it as a sum of sines and cosines. That's exactly what the discrete-time Fourier transform does:

$$X_{2\pi}(\omega) = \sum_{n=-\infty}^{\infty} x[n]e^{-i\omega n}$$

Note that this is just a special case of the Z-transform, where $z = e^{i\omega}$.

The frequency response of the filter describes how the spectrum of a signal is altered when it passes through the filter. It relates the spectrum of the output signal $Y(\omega)$ to the spectrum of the input signal $X(\omega)$. We already had an expression for the spectrum of the output divided by the spectrum of the input, we called it the transfer function $H(z) = Y(z)/X(z)$. To get the frequency response of the filter, we can just evaluate the transfer function for $z = e^{i\omega}$. Also note that this is the DTFT of the impulse response $h[k]$.

We use the formula $e^{j\theta} - e^{-j\theta} = 2j \sin(\theta)$.

$$\begin{aligned} \mathcal{F}_{DTFT}\{h[k]\} &= \sum_{n=-\infty}^{\infty} h[n]e^{-i\omega n} \\ &= H(e^{i\omega}) \\ &= \frac{1}{N} \frac{1 - e^{-j\omega N}}{1 - e^{-j\omega}} \\ &= \frac{1}{N} \frac{e^{-j\omega N/2}}{e^{-j\omega/2}} \frac{e^{j\omega N/2} - e^{-j\omega N/2}}{e^{j\omega/2} - e^{-j\omega/2}} \\ &= \frac{1}{N} \frac{e^{-j\omega N/2}}{e^{-j\omega/2}} \frac{2j \sin\left(\frac{\omega N}{2}\right)}{2j \sin\left(\frac{\omega}{2}\right)} \\ &= \frac{e^{-j\frac{\omega}{2}(N-1)}}{N} \frac{\sin\left(\frac{\omega N}{2}\right)}{\sin\left(\frac{\omega}{2}\right)} \end{aligned}$$

We can now calculate the amplitude of each frequency component in the output by taking the modulus of $H(e^{i\omega})$. For reasons that will become apparent in a minute, we'll calculate the modulus squared.

$$\begin{aligned} |H(e^{i\omega})|^2 &= \left| \frac{e^{-j\frac{\omega}{2}(N-1)}}{N} \frac{\sin\left(\frac{\omega N}{2}\right)}{\sin\left(\frac{\omega}{2}\right)} \right|^2 \\ &= \frac{1}{N^2} \frac{\sin^2\left(\frac{\omega N}{2}\right)}{\sin^2\left(\frac{\omega}{2}\right)} \end{aligned}$$

ω is the normalized pulsantance (angular frequency) in radians per sample. You can substitute it with $\omega = \frac{2\pi f}{f_s}$ where f is the frequency in Hertz, and f_s is the sample frequency of the system in Hertz.

We can now plot the filter's gain in function of the frequency. These plots often use a logarithmic scale, to show the gain in decibels. In order to calculate the power gain, the amplitude is squared.

$$A_{dB}(\omega) = 10 \log_{10} |H(e^{i\omega})|^2$$

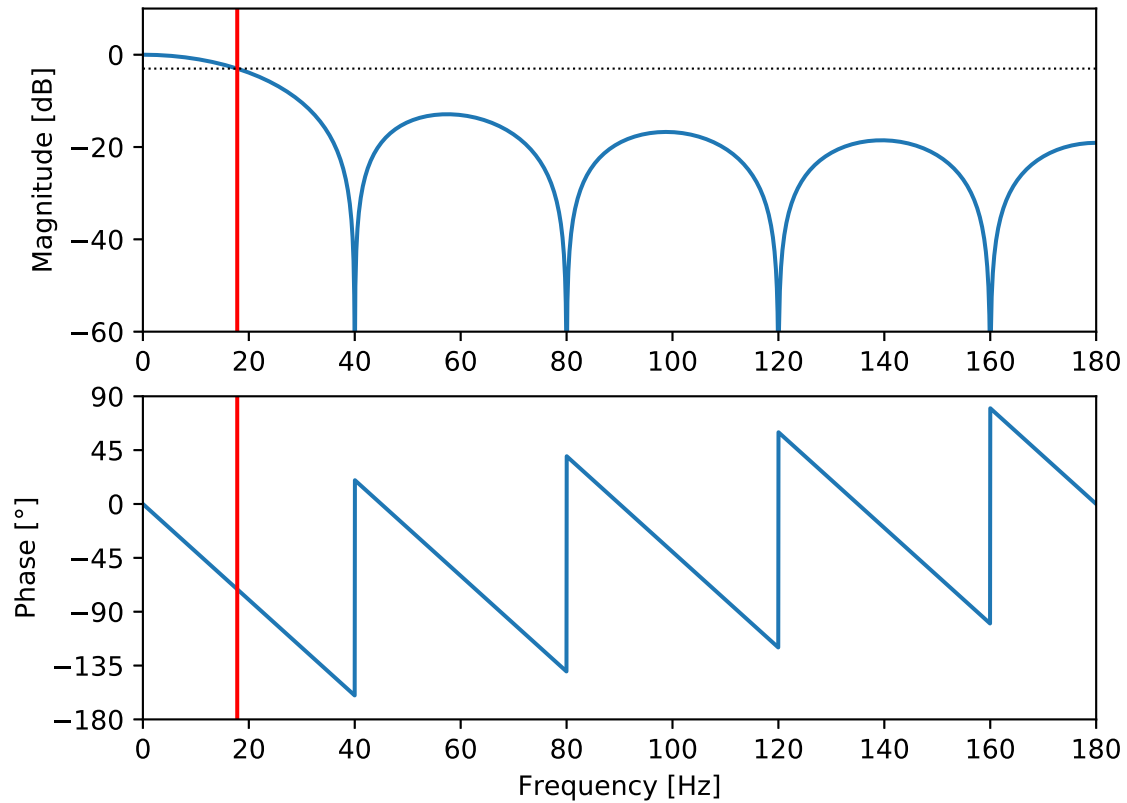
Note that when a frequency is not present in the output signal, the gain will be $-\infty$ dB. If a frequency has an amplitude of one in the output signal, the gain will be 0 dB.

We can also plot the phase shift introduced by the SMA. This is just the complex argument of the transfer function.

$$\phi(\omega) = \angle H(e^{i\omega})$$

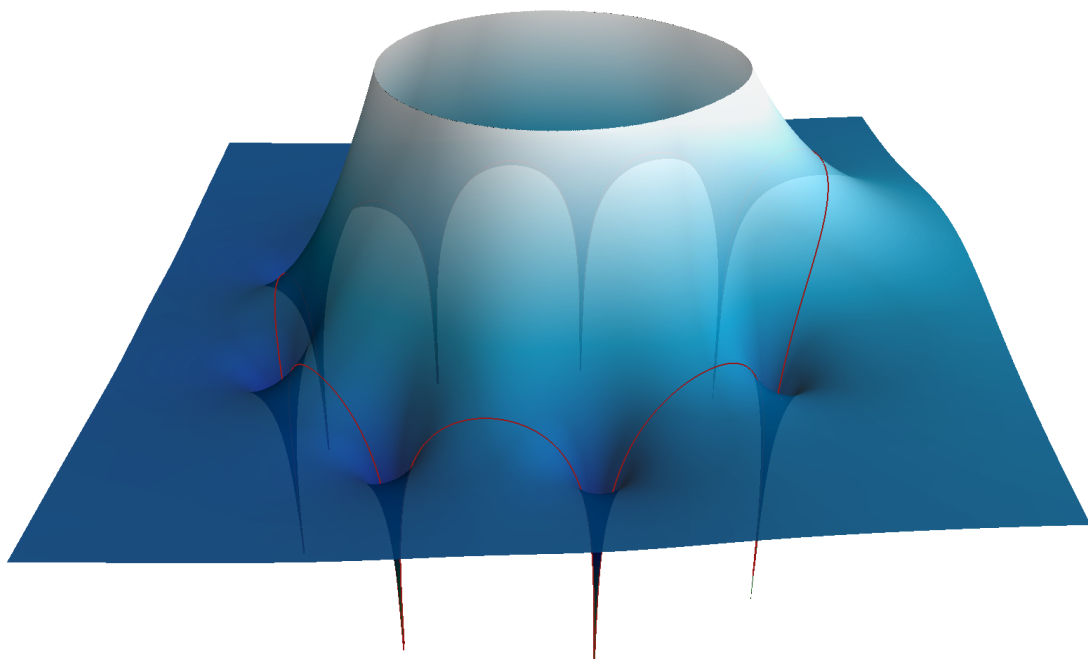
The image below shows the bode plot of an SMA with a length of $N = 9$, and a sample frequency of $f_s = 360\text{Hz}$. The frequency axis is limited between 0 Hz and the Nyquist frequency $f_s/2$. This corresponds to a normalized pulsance of $\omega \in [0, \pi)$.

Bode Plot



You can clearly see the low-pass behavior of the SMA: low frequencies have a near-unit gain, and high frequencies are attenuated. Also note the zeros, where the magnitude plot goes to $-\infty\text{ dB}$. In a zero, the transfer function changes sign, so the phase jumps 180° as well.

To get a better understanding of where this curve comes from, we can plot the entire $|H(z)|$ surface in the Z-domain. As mentioned above, the DTFT is just a special case of the Z-transform, where $z = e^{i\omega}$, i.e. the unit circle in the complex plane. Remember that the point $z = e^{i\omega}$ is a point with a distance of 1 to the origin, and with an angle of $\omega\text{ rad}$ between its position vector and the positive x axis.



[Image source code](#)

The image of the unit circle is shown in red. Notice that this is the same curve as the blue curve in the magnitude response graph above: close to 0 when $\omega \rightarrow 0$ (the right half of the circle) and negative when $\omega \rightarrow \pi$ (the left half of the circle). The positive x axis points to the right-hand side, the y axis points away from you, and the z axis points up).

Cutoff frequency

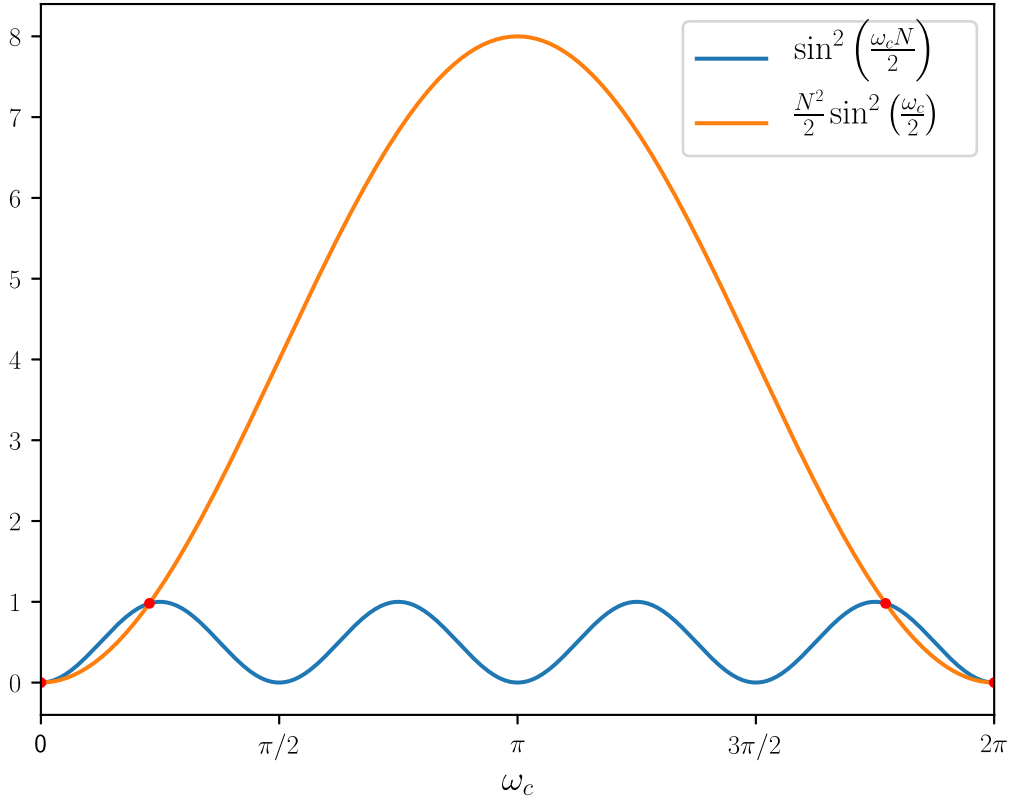
The cutoff frequency is defined as the frequency of the half-power point, where the power gain is a half. It's often called the -3 dB -point, because $10 \log_{10} \left(\frac{1}{2} \right) \approx -3.01 \text{ dB}$.

To find it, just solve the following equation:

$$\begin{aligned} |H(e^{i\omega_c})|^2 &= \frac{1}{2} \\ \frac{1}{N^2} \frac{\sin^2 \left(\frac{\omega_c N}{2} \right)}{\sin^2 \left(\frac{\omega_c}{2} \right)} &= \frac{1}{2} \\ \sin^2 \left(\frac{\omega_c N}{2} \right) &= \frac{N^2}{2} \sin^2 \left(\frac{\omega_c}{2} \right) \\ \sin^2 \left(\frac{\omega_c N}{2} \right) - \frac{N^2}{2} \sin^2 \left(\frac{\omega_c}{2} \right) &= 0 \end{aligned}$$

This equation doesn't have a general analytical solution, but we can use numerical methods to determine the cut-off frequency.

Solutions of $\sin^2 \left(\frac{\omega_c N}{2} \right) = \frac{N^2}{2} \sin^2 \left(\frac{\omega_c}{2} \right)$ for $N = 4$



Note that the intersections at zero and 2π are not valid solutions. In these points the transfer function is undefined, because it evaluates to $0/0$. We only care about the first solution for $\omega_c \in (0, \pi)$.

A possible approach to find the intersection is to use the [Newton-Raphson method](#).

A good starting point seems to be at half the period of the high-frequency sine wave.

We don't have to keep the squares, because both sines will be positive in the region of the intersection.

$$\begin{aligned} f(\omega_c) &= \sin \left(\frac{\omega_c N}{2} \right) - \frac{N}{\sqrt{2}} \sin \left(\frac{\omega_c}{2} \right) \\ f'(\omega_c) &= \frac{N}{2} \cos \left(\frac{\omega_c N}{2} \right) - \frac{N}{2\sqrt{2}} \cos \left(\frac{\omega_c}{2} \right) \\ \omega_c^{(0)} &= \pi/N \end{aligned}$$

A Python implementation can be found in the `get_sma_cutoff` function in the snippet below.

Plotting the frequency response, impulse response and step response in Python

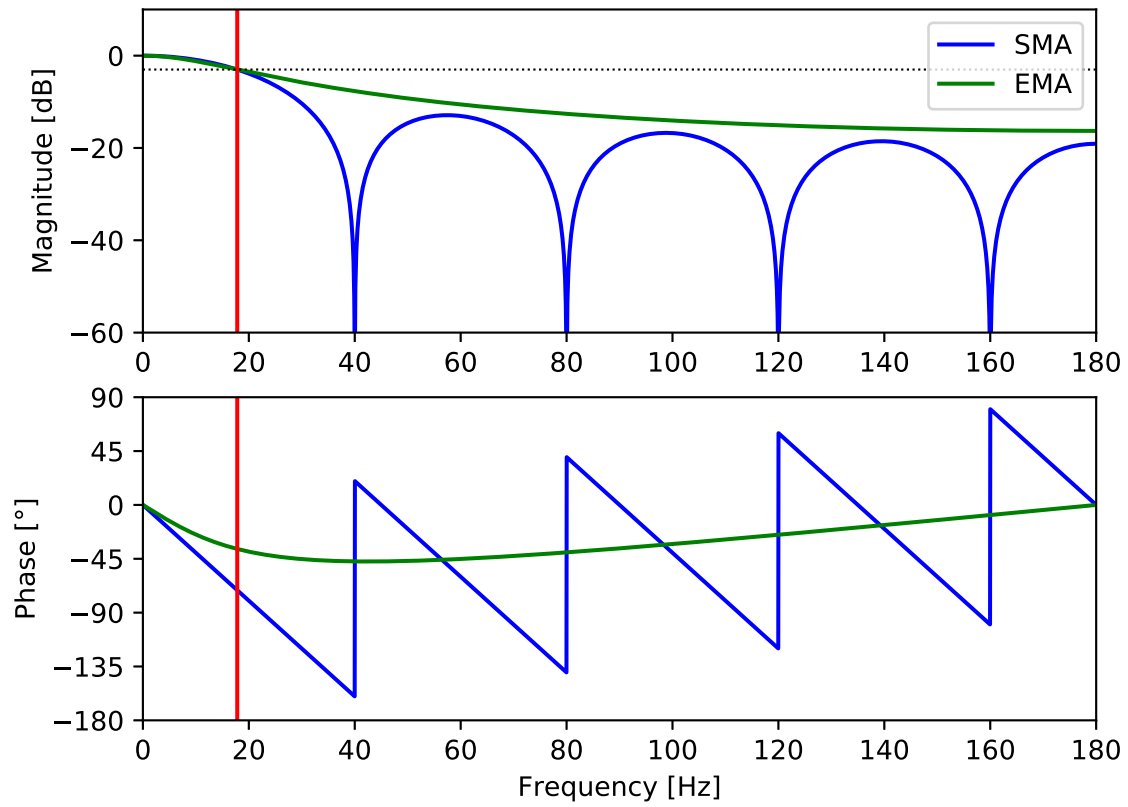
We can use the SciPy and Matplotlib modules to plot the frequency response in Python. This is the script that was used to generate the plots in the previous paragraphs.

```
1 from scipy.optimize import newton
2 from scipy.signal import freqz, dimpulse, dstep
3 from math import sin, cos, sqrt, pi
4 import numpy as np
5 import matplotlib.pyplot as plt
6
7 # Function for calculating the cut-off frequency of a moving average filter
8 def get_sma_cutoff(N, **kwargs):
9     func = lambda w: sin(N*w/2) - N/sqrt(2) * sin(w/2) #  $|H(e^{j\omega})| = \sqrt{2}/2$ 
10    deriv = lambda w: cos(N*w/2) * N/2 - N/sqrt(2) * cos(w/2) / 2 #  $dfunc/dx$ 
11    omega_0 = pi/N # Starting condition: halfway the first period of  $\sin(N\omega/2)$ 
12    return newton(func, omega_0, deriv, **kwargs)
13
14 # Simple moving average design parameters
15 f_s = 100
16 N = 5
17
18 # Find the cut-off frequency of the SMA
19 w_c = get_sma_cutoff(N)
20 f_c = w_c * f_s / (2 * pi)
21
22 # SMA coefficients
23 b = np.ones(N)
24 a = np.array([N] + [0]*(N-1))
25
26 # Calculate the frequency response
27 w, h = freqz(b, a, worN=4096)
28 w *= f_s / (2 * pi) # Convert from rad/sample to Hz
29
30 # Plot the amplitude response
31 plt.subplot(2, 1, 1)
32 plt.suptitle('Bode Plot')
33 plt.plot(w, 20 * np.log10(abs(h))) # Convert modulus to dB
34 plt.ylabel('Magnitude [dB]')
35 plt.xlim(0, f_s / 2)
36 plt.ylim(-60, 10)
37 plt.axvline(f_c, color='red')
38 plt.axhline(-3.01, linewidth=0.8, color='black', linestyle=':')
39
40 # Plot the phase response
41 plt.subplot(2, 1, 2)
42 plt.plot(w, 180 * np.angle(h) / pi) # Convert argument to degrees
43 plt.xlabel('Frequency [Hz]')
44 plt.ylabel('Phase [°]')
45 plt.xlim(0, f_s / 2)
46 plt.ylim(-180, 90)
47 plt.yticks([-180, -135, -90, -45, 0, 45, 90])
48 plt.axvline(f_c, color='red')
49 plt.show()
50
51 # Plot the impulse response
52 t, y = dimpulse((b, a, 1/f_s), n=2*N)
53 plt.suptitle('Impulse Response')
54 _, _, baseline = plt.stem(t, y[0], basefmt='k:')
55 plt.setp(baseline, 'linewidth', 1)
56 baseline.set_xdata([0,1])
57 baseline.set_transform(plt.gca().get_yaxis_transform())
58 plt.xlabel('Time [seconds]')
59 plt.ylabel('Output')
60 plt.xlim(-1/f_s, 2*N/f_s)
61 plt.yticks([0, 0.5/N, 1.0/N])
62 plt.show()
63
64 # Plot the step response
65 t, y = dstep((b, a, 1/f_s), n=2*N)
66 plt.suptitle('Step Response')
67 _, _, baseline = plt.stem(t, y[0], basefmt='k:')
68 plt.setp(baseline, 'linewidth', 1)
69 baseline.set_xdata([0,1])
70 baseline.set_transform(plt.gca().get_yaxis_transform())
71 plt.xlabel('Time [seconds]')
72 plt.ylabel('Output')
73 plt.xlim(-1/f_s, 2*N/f_s)
74 plt.yticks([0, 0.2, 0.4, 0.6, 0.8, 1])
75 plt.show()
```

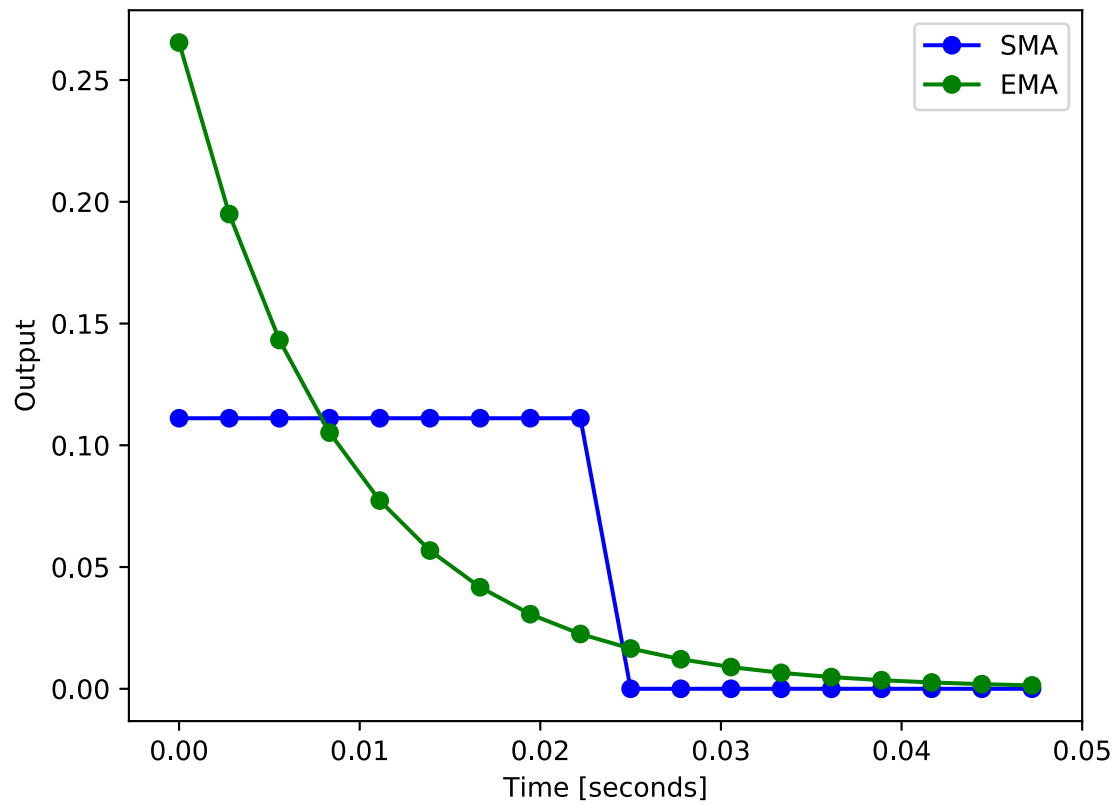
Comparing the Simple Moving Average filter to the Exponential Moving Average filter

Using the same Python functions as before, we can plot the responses of the EMA and the SMA on top of each other. First, the length N of the SMA is chosen, then its 3 dB cut-off frequency is calculated, and this frequency is then used to design the EMA. Do note that this is a fairly arbitrary decision.

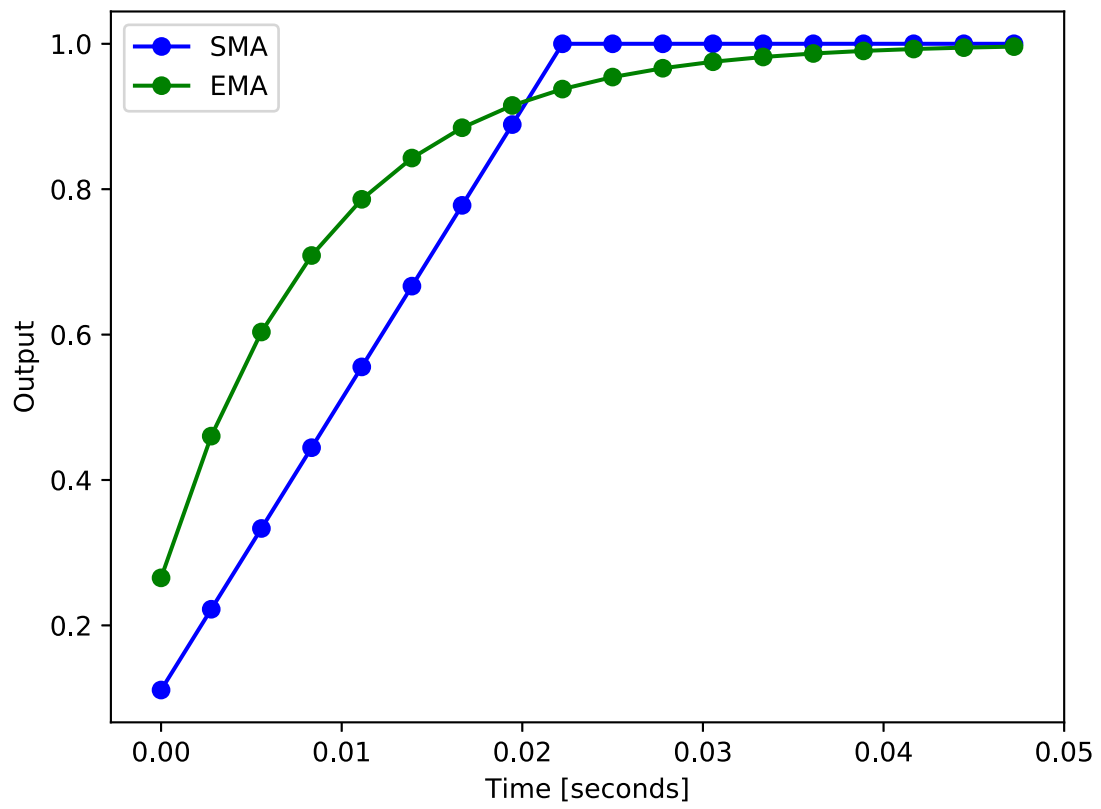
Bode Plot



Impulse Response



Step Response



```

1  from scipy.optimize import newton
2  from scipy.signal import freqz, dimpulse, dstep
3  from math import sin, cos, sqrt, pi
4  import numpy as np
5  import matplotlib.pyplot as plt
6
7  # Function for calculating the cut-off frequency of a moving average filter
8  def get_sma_cutoff(N, **kwargs):
9      func = lambda w: sin(N*w/2) - N/sqrt(2) * sin(w/2) #  $|H(e^{j\omega})| = \sqrt{2}/2$ 
10     deriv = lambda w: cos(N*w/2) * N/2 - N/sqrt(2) * cos(w/2) / 2 #  $d\text{func}/dw$ 
11     omega_0 = pi/N # Starting condition: halfway the first period of  $\sin(N\omega/2)$ 
12     return newton(func, omega_0, deriv, **kwargs)
13
14 # Simple moving average design parameters
15 f_s = 360
16 N = 9
17
18 # Find the cut-off frequency of the SMA
19 w_c = get_sma_cutoff(N)
20 f_c = w_c * f_s / (2 * pi)
21
22 # Calculate the pole location of the EMA with the same cut-off frequency
23 b = 2 - 2*cos(w_c)
24 alpha = (-b + sqrt(b**2 + 4*b)) / 2
25
26 # SMA & EMA coefficients
27 b_sma = np.ones(N)
28 a_sma = np.array([N] + [0]*(N-1))
29
30 b_ema = np.array((alpha, 0))
31 a_ema = np.array((1, alpha - 1))
32
33 # Calculate the frequency response
34 w, h_sma = freqz(b_sma, a_sma, worN=4096)
35 w, h_ema = freqz(b_ema, a_ema, w)
36 w *= f_s / (2 * pi) # Convert from rad/sample to Hz
37
38 # Plot the amplitude response
39 plt.subplot(2, 1, 1)
40 plt.suptitle('Bode Plot')
41 plt.plot(w, 20 * np.log10(abs(h_sma)), # Convert modulus to dB
42         color='blue', label='SMA')
43 plt.plot(w, 20 * np.log10(abs(h_ema)),
44         color='green', label='EMA')
45 plt.ylabel('Magnitude [dB]')
46 plt.xlim(0, f_s / 2)
47 plt.ylim(-60, 10)
48 plt.axvline(f_c, color='red')
49 plt.axhline(-3.01, linewidth=0.8, color='black', linestyle=':')
50 plt.legend()
51
52 # Plot the phase response
53 plt.subplot(2, 1, 2)
54 plt.plot(w, 180 * np.angle(h_sma) / pi, # Convert argument to degrees
55         color='blue')
56 plt.plot(w, 180 * np.angle(h_ema) / pi,
57         color='green')
58 plt.xlabel('Frequency [Hz]')
59 plt.ylabel('Phase [°]')
60 plt.xlim(0, f_s / 2)
61 plt.ylim(-180, 90)
62 plt.yticks([-180, -135, -90, -45, 0, 45, 90])
63 plt.axvline(f_c, color='red')
64 plt.show()
65
66 # Plot the impulse response
67 t, y_sma = dimpulse((b_sma, a_sma, 1/f_s), n=2*N)
68 t, y_ema = dimpulse((b_ema, a_ema, 1/f_s), n=2*N)
69 plt.suptitle('Impulse Response')
70 plt.plot(t, y_sma[0], 'o-',
71         color='blue', label='SMA')
72 plt.plot(t, y_ema[0], 'o-',
73         color='green', label='EMA')
74 plt.xlabel('Time [seconds]')
75 plt.ylabel('Output')
76 plt.xlim(-1/f_s, 2*N/f_s)
77 plt.legend()
78 plt.show()
79
80 # Plot the step response
81 t, y_sma = dstep((b_sma, a_sma, 1/f_s), n=2*N)
82 t, y_ema = dstep((b_ema, a_ema, 1/f_s), n=2*N)
83 plt.suptitle('Step Response')
84 plt.plot(t, y_sma[0], 'o-',
85         color='blue', label='SMA')
86 plt.plot(t, y_ema[0], 'o-',
87         color='green', label='EMA')
88 plt.xlabel('Time [seconds]')
89 plt.ylabel('Output')
90 plt.xlim(-1/f_s, 2*N/f_s)
91 plt.legend()
92 plt.show()

```