

COE3DY4 Lab #1

DSP Primitives in Python

Objective

The purpose of this lab is to develop the understanding of some basic digital signal processing (DSP) primitives used in Software Defined Radios (SDRs) and implement them in Python.

Preparation

- Revise the fundamental material from signals and systems, in particular, the linear time-invariant (LTI) systems, impulse response, convolution, Nyquist rate, Fourier transform
- Revise basic programming in both Python and MATLAB, because of the inherent similarities between NumPy/SciPy and MATLAB
- Get familiar with NumPy (<https://numpy.org/>) at SciPy (<https://scipy.org/>)
- For visualizing your data in Python refer to the documentation for Matplotlib (<https://matplotlib.org/>)

Discrete Fourier Transform

The Fourier transform is widely used in science and engineering, from physics and astronomy to multimedia processing and communications to data science and machine learning. One simple way to think of the Fourier transform is as a change of a signal representation from one domain (time) to another one (frequency) for more intuitive data visualization and/or more efficient computational processing. It relies on complex exponentials, used to manipulate periodic trigonometric functions, i.e., $\cos(\theta)$ and $\sin(\theta)$, using Euler's formula $e^{j\theta} = \cos(\theta) + j \sin(\theta)$, where $j = \sqrt{-1}$ is the imaginary unit.

Formally, the continuous-time Fourier Transform of a time domain signal $x(t)$ is defined as:

$$X(f) = \int_{-\infty}^{\infty} x(t) e^{-j2\pi ft} dt \quad \text{or} \quad X(\omega) = \int_{-\infty}^{\infty} x(t) e^{-j\omega t} dt$$

$X(f)$ or $X(\omega)$ gives the signal's frequency spectrum, where the frequency f is measured in Hz, and the angular frequency ω is measured in radians/second. The above is also called the *analysis* formula. The *synthesis* formula is defined using the inverse Fourier transform:

$$x(t) = \int_{-\infty}^{\infty} X(f) e^{j2\pi ft} df \quad \text{or} \quad x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(\omega) e^{j\omega t} d\omega$$

In Fourier analysis, the time and frequency domain variables can be continuous or discrete. The above formulas assume continuous time/frequency across infinite intervals. A detailed and rigorous

mathematical coverage of the Fourier transform, which goes beyond the scope of this project course on SDRs, can be found online at <https://ccrma.stanford.edu/~jos/st/>. For our project, it is still worth stating that the discrete-time Fourier Transform discretizes the time variable only, whereas the Discrete Fourier Transform (DFT) has both the time and the frequency variables discretized. Since DFT is the most commonly used Fourier-type transform in practical applications due to its computational nature, we will use it for the rest of this project on SDRs.

Assuming we have an input signal represented as a block of N samples $x(k)$, with $k = 0 \dots N - 1$, the DFT is defined as:

$$X(m) = \sum_{k=0}^{N-1} x(k) e^{-2\pi j \frac{km}{N}}$$

where $X(m)$ is the m^{th} frequency component (or bin) in the spectrum of the input signal. The above formula has a straightforward translation to pseudocode:

Input: Time domain samples $x(k)$ with $k = 0 \dots N - 1$

Output: Discrete frequency bins $X(m)$ with $m = 0 \dots N - 1$

```

for  $m \in [0, N - 1]$  do
     $X(m) \Leftarrow 0$ 
    for  $k \in [0, N - 1]$  do
         $X(m) \Leftarrow X(m) + x(k) \times e^{-2\pi j \frac{km}{N}}$ 
    end for
end for

```

Note, using the notation of **twiddle factors**, i.e., $W_N^r = e^{2\pi j \frac{r}{N}}$, the DFT can be rewritten as:

$$X(m) = \sum_{k=0}^{N-1} x(k) W_N^{-km}$$

W_N^r are in fact the N^{th} roots of unity and simple math can show that $(W_N^r)^N = 1$, for each $r = 0 \dots N - 1$. Note, the properties of these twiddle factors are central to the derivation of the faster implementations of the DFT, such as the Fast Fourier Transform (FFT) ¹.

Based on the same inversion principle as for the continuous-time Fourier transform, the Inverse Discrete Fourier Transform (IDFT) is defined as:

$$x(k) = \frac{1}{N} \sum_{m=0}^{N-1} X(m) W_N^{km}$$

In the lab, you should perform the following tasks, which are aimed at strengthening the link between *mathematical formalism* and *algorithmic thinking* in DSP (and system design in general). The command line arguments are documented in the main function in `fourierTransform.py` and they are: reference code (*rc*), in-lab experiment 1 (*il1*), in-lab experiment 2 (*il2*) and in-lab experiment 3 (*il3*).

- (*rc*): Run the reference code using the appropriate command line argument (*rc*) and develop an understanding of how sine waves are plotted in both time and frequency domains.

¹Although **twiddle factors** can be defined also as $W_N^r = e^{-2\pi j \frac{r}{N}}$ the same principles for deriving FFT hold.

- (*il1*): Write your **own** functions for DFT and IDFT in Python. Confirm their functional correctness against the built-in implementations of the FFT and inverse FFT (IFFT) from SciPy by monitoring on the screen the data in both time and frequency domains.
- (*il2*): Generate a random signal using NumPy with 1,000 samples, whose magnitude is normalized to -10 to +10, and confirm that in your implementation of DFT/IDFT, the signal energy holds (as expected from Parseval's theorem). Note that the signal energy in the time domain is the summation of the squares of samples' absolute values, i.e., $|x(k)|^2$. In the frequency domain, it is the summation across all frequency components of their individual energy spectral density, i.e., $\frac{1}{N} |X(m)|^2$. Stated differently, $\sum_{k=0}^{N-1} |x(k)|^2$ must equal $\frac{1}{N} \sum_{m=0}^{N-1} |X(m)|^2$.
- (*il3*): Using the reference code for a pure tone (a sine wave) that is plotted in time and frequency domains, write a function to generate a multi-tone signal (up to 3 tones is sufficient) where the sines have different amplitudes, frequencies and phases. Visualize your multi-tone signal in both time and frequency domains.

Digital Filter Design

For linear-time invariant (LTI) systems that are stable and causal, digital filtering can be performed via convolution (using operator $*$). Consider two signals $x(a)$ of size A and $h(b)$ of size B . The discrete convolution of these two sequences $y(c)$ is defined as :

$$y(c) = (x * h)(c) = \sum_{l=-\infty}^{\infty} x(l)h(c-l) = \sum_{l=-\infty}^{\infty} h(l)x(c-l) = (h * x)(c)$$

In the above formula, when the indices for x and y are below zero or above their maximum size, i.e., $A-1$ for x and $B-1$ for h , the values of the corresponding $x(l)$ and $h(c-l)$ are assumed to be zero. The total number of non-zero elements in signal y will equal $A+B-1$.

A central theorem of the Fourier theory is that the transform of the convolution of two sequences in the time domain is equal to the pointwise multiplication of their transforms in the frequency domain (referred to as the convolution theorem). One of the implications of this theorem is:

$$DFT(x * h) = DFT(x)DFT(h) \implies x * h = IDFT(DFT(x)DFT(h))$$

The above can be leveraged for the computational speed-up of the convolution operation (only when DFT is replaced via FFT). However, since computational speed-ups are beyond this lab's scope, the convolution theorem will be used to understand how to derive the impulse response of a digital filter, i.e., the coefficients that are used to weigh the input samples.

An ideal low-pass filter will eliminate all the input signal frequency components above a cutoff frequency f_c . This can be achieved by multiplying the frequency response of the input sequence with a rectangular window; the height of the rectangular window defines the filter gain, and the window width from $-f_c$ to f_c defines the filter's pass band. The inverse Fourier transform of the rectangular window is the *sinc* function, defined as $\frac{\sin(\pi t)}{\pi t}$ (when t is 0 *sinc* is defined by its limiting value of 1). The shapes of the rectangular window and the *sinc* function are shown on the next page for illustrative purposes.

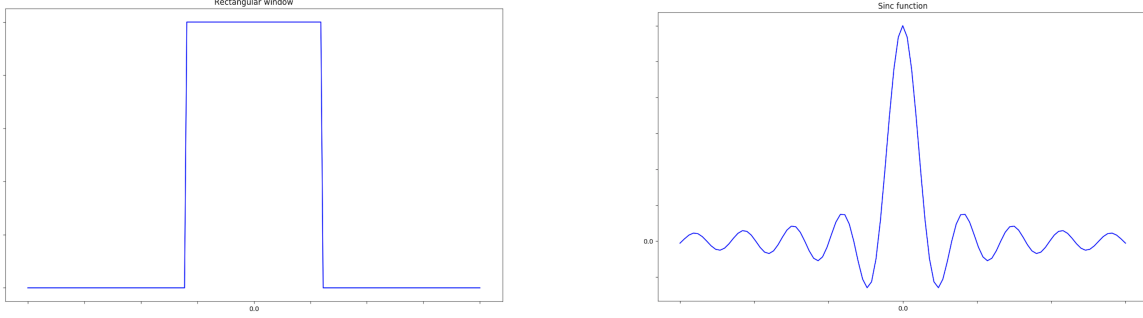


Figure 1: The Fourier transform pair: rectangular window and *sinc* function

Based on the above, the impulse response of a low-pass filter is derived from the *sinc* function, whose argument and scale factor are dependent on the normalized cutoff frequency as follows:

$$h(i) = \frac{f_c}{f_s/2} \text{sinc}\left(\frac{f_c}{f_s/2}i\right)$$

In the above formula, f_c is the cutoff frequency, and f_s is the sampling rate. Note, when computing the normalized frequency $\frac{f_c}{f_s/2}$ the cutoff frequency f_c is divided by $f_s/2$, which is the Nyquist frequency that must be greater than the bandwidth of the input signal. Note also, from the above formula, it is implied the impulse response converges to $\frac{f_c}{f_s/2}$ at its center point. Note also that i is only an index in a vector and has nothing to do with the imaginary numbers.

One point that must be clarified before translating the previous mathematical formalism into pseudocode is that the impulse response derived via the *sinc* function needs to be *windowed*. In practice, the filters based on convolution are limited in size, hence the term finite impulse response (FIR) filters. Because the ideal filter assumes an infinite response, FIRs will truncate the *sinc* function based on the number of filter's taps N_{taps} . If a rectangular window is used for truncation, one of its well-understood side-effects is excessive *spectral leakage* in the frequency domain. Spectral leakage is the phenomenon where new frequencies arise in the spectrum because of the behaviour induced by truncation at the boundaries of the FIR. To alleviate this problem, the impulse response is *windowed*. The trade-offs involved in different types of windows and their impact on spectral leakage are well-studied, and in our filter design, we will rely on the Hann window, one of the oldest yet most commonly used in practice.

The pseudocode for deriving the impulse response coefficients of a low pass filter is given below.

Input: Filter parameters: f_c , f_s and N_{taps}

Output: Taps of the FIR filter $h(i)$ with $i = 0 \dots N_{taps} - 1$

$Norm_{cutoff} \leftarrow \frac{f_c}{f_s/2}$ ▷ define explicitly the normalized cutoff frequency $Norm_{cutoff}$

for $i \in [0, N_{taps} - 1]$ **do**

if $i = (N_{taps} - 1)/2$ **then**

$h(i) \leftarrow Norm_{cutoff}$ ▷ avoid division by zero in *sinc* for the center tap when N_{taps} is odd

else

$h[i] \leftarrow Norm_{cutoff} \frac{\sin(\pi Norm_{cutoff}(i - (N_{taps} - 1)/2))}{\pi Norm_{cutoff}(i - (N_{taps} - 1)/2)}$

end if

$h(i) \leftarrow h(i) \sin^2\left(\frac{i\pi}{N_{taps}}\right)$

▷ apply the Hann window

end for

It is important to mention that deciding on the number of taps for an FIR filter can be a cumbersome task. It depends on many factors, including passband gain/ripple, stopband attenuation/ripple, and transition band width. On the one hand, the larger the number of taps, the better the FIR approximation of the infinite impulse response, hence better signal quality; on the other hand, the more taps, the more multiplications and accumulations need to be done, hence there is an increased computational demand. An empirical yet effective decision can be reached using the "harris rule of thumb", where the number of taps is determined as:

$$N_{taps} = \frac{f_s}{\Delta f_t} \frac{Attn(dB)}{22}$$

In the above formula, in addition to the already-introduced number of taps (N_{taps}) and the sampling frequency (f_s), the Δf_t is the transition band width and $Attn(dB)$ is the stop band attenuation (in dB from the main lobe to the peak of the first side lobe in the magnitude frequency response). For our SDR application, usually at least $40dB$ attenuation is preferred; naturally, the higher the attenuation the better it is in terms of signal quality, however it is worse in terms of compute speed, which makes it more challenging to meet real-time constraints. The relative size of the transition band width to the sampling rate varies from one use case to another, however it is rarely the case that Δf_t is much higher than 10% of f_s . To conclude, the key point of this discussion is that, in practice, compromises are needed when implementing filters in real-time by trading off the signal quality against the compute speed. Other non-functional requirements, e.g., energy-efficiency, can add further dimensions to this complex engineering problem of system design.

Consistent with the earlier objective to link *mathematical formalism* to *algorithmic thinking*, in the lab, you should modify `filterDesign.py` to perform the following.

- (*rc*): Run the reference code and visualize the frequency response of the low-pass filter designed through `firwin`. Start varying the sample rate, the cutoff frequency and the number of filter taps, and understand the changes in the shape of the frequency response for the low-pass filter (in particular, the width of the transition band and attenuation)
- (*il1*): Implement the given pseudocode to design your **own** low-pass impulse response in **Python**. Compare your impulse response against the built-in methods from **SciPy** for filter design, such as `firwin`. This can be achieved by visually inspecting the frequency response of both methods using `freqz`.
- (*il2*): Use your impulse response to filter out the highest frequency of the multi-tone signal from the previous experiment. You can achieve this by using `lfilter`. You should visualize the outcome both in the time and frequency domains.

Digital Filtering of Data Streams Divided into Blocks

The previous sections have discussed a few foundational concepts in DSP, which have enabled you to build your own tools for the design of impulse responses and spectrum analysis. They will be useful beyond this lab as you transition to a more refined implementation in C/C++. This section discusses a complementary yet important practical consideration that must be accounted for in real-time implementations when input data is continuously streamed.

The impulse response convolution with an input data sequence can be done in a single-pass only if **all** the input data has been buffered. This is not practically feasible in many applications,

especially when the input is streamed. In such cases, the input stream is segmented into **blocks** that are processed one at a time. However, to avoid discontinuities in the output sequence, at the end of each block, the convolution needs to be stopped and resumed from a state that contains all the input samples needed to compute the outputs for the new block.

The above challenge is best explained through an example. Assume the input x is convolved with an impulse response h with 10 taps. Assume also the blocks are of size 100. Consider the case when two consecutive output samples from different blocks ($y(99)$ and $y(100)$) are computed:

$$y(99) = x(90)h(9) + x(91)h(8) + \dots + x(99)h(0)$$

$$y(100) = x(91)h(9) + x(92)h(8) + \dots + x(100)h(0)$$

For $y(99)$, all the input samples are from the same block. However, for $y(100)$ we need one sample from the second block ($x(100)$) and 9 samples from the first block ($x(91)$ to $x(99)$). If the convolution for the second block is processed independently from the first one, then samples $x(91)$ to $x(99)$ will be missing (they are assumed to be zero by default), which clearly introduces discontinuities in the output sequence. To avoid the above problem, there is the need to save the samples $x(91)$ to $x(99)$ from the first block to use them when computing the first sample from the second block.

In the lab, you should modify `blockProcessing.py` to perform the following:

- (*rc*): The reference code filters an audio `.wav` file in a single-pass. Understand how the entire code is set up (including where the `.wav` files are read and written). You can use an audio player to listen to them. Vary the cutoff frequency for the low-pass filter and understand the impact on the sound quality.
- (*il1*) The data is processed in blocks, where the filter state at the end of each block is saved only for one of the two audio channels. State saving is achieved by using the `zi` argument for `lfilter` (check SciPy documentation for further details). Using the given code, you can remove and add state saving to the experiment as follows: vary the size of the blocks and understand when and why there are audible artifacts due to a lack of state saving when transitioning between consecutive blocks. When state saving is performed, confirm that the output audio file produces the same sound quality as when using single-pass filtering. Note that it is fair to assume that the block size will always be greater than the filter size (as it is common in practice).

Take-home exercises

In addition to **completing** and **submitting** your in-lab experiments, to further improve your understanding of DSP, as well as to become more proficient with **Python** and comfortable with its libraries for signal processing, you should perform the following:

- Add functionality to `fourierTransform.py` to generate periodic square waves of arbitrary duty cycles. The arguments to this `generateSquare` function should be the same ones as for the function from the reference code used for generating the sine wave (`generateSin`), with the exception that the phase argument is replaced with a duty cycle argument, which takes a real value greater than zero and smaller than one (as a default use a duty cycle of 0.5). Take note that for this exercise, you are **allowed** to use functions from **NumPy** and **SciPy** to implement the square wave faster; as a suggestion, explore the `square` function from the

`scipy.signal` library. Plot these waves in both the time and frequency domains and try to understand what is going on when you vary the arguments of your function, as well as the default values for the sampling rate. As you try to explain the frequency spectrum of square waves, revisit your background knowledge from signal theory, search for explanations, consolidate your understanding, ..., and summarize your thoughts in the report.

- Learn how to implement a bandpass filter using `firwin` from `SciPy`. Confirm your understanding by plotting the frequency response in `filterDesign.py` and looking at filtered data in the time domain. To validate your understanding, you must pass only the second lowest tone from a four-tone signal (by removing the lowest and the two highest frequency tones). As you explore, perform the same type of experiments as for the in-lab experiments, i.e., vary the number of filter taps, sampling rate and cutoff frequencies (i.e., get the cutoff frequencies closer to and further away from the passband edges). Summarize your thoughts in the report by relating your background knowledge to the exploratory work from this undertaking.
- Start by performing audio filtering in `blockProcessing.py` by replacing `firwin` with your own function (from the in-lab experiment on filter design) to derive the filter coefficients. Furthermore, you should implement your own function for convolution that should replace the functionality from `lfilter`, as well as it will guarantee that data can be processed in blocks without discontinuities in the output sequence (assume the block size is always greater than the number of filter coefficients). It is not expected for your solution to derive the filter coefficients and convolution to produce bit-level equivalence against `firwin` and `lfilter`. Note also that the computational speed is of no concern (at this time in the course). Nevertheless, the filtered audio must be perceived in the same way when using your own methods vs `firwin` and `lfilter` under the same filter and convolution parameters (i.e., number of taps, cutoff frequency and block size). Summarize in the report the challenges and your solutions for this (more complex) development task.

It is expected to build your own library of DSP primitives, i.e., discrete Fourier transform, impulse response coefficients and block convolution, incrementally with regular self-checks and comparisons against the reference functions from `SciPy`. This is critical to smoothen the transition from modelling DSP primitives in `Python` to their implementation in C/C++ for an SDR system, which will be the focus of the next lab.

The report must have three sections, one for each of the above items. One page is sufficient and should not be exceeded unless there are out-of-ordinary points to be made.

The source code should stay in the `model` sub-folder of the GitHub repo; your report (in `.pdf`, `.txt` or `.md` format) should be included in the `doc` sub-folder. The source code already provides the expected command-line arguments that should be used for in-lab experiments/take-home exercises.

Your submission, which is to be done via a push to the master branch of your group's GitHub repo for this lab, is due 16 hours before your next lab session is scheduled. Late submissions will be penalized.

This lab is worth 6 % of your total grade.