



University
of Glasgow | School of
Engineering

Digital Signal Processing

Bernd Porr & Nick Bailey

October 28, 2021

Creative commons BY-SA (C) Bernd Porr, bernd.porr@glasgow.ac.uk

Contents

1	Acknowledgements	4
2	Introduction	4
2.1	Suggested reading	4
2.2	Advantages of digital signal processing	4
2.3	Development boards	4
3	Python	4
3.1	How to use Python?	5
3.2	Packages/Modules	6
3.3	Data types	7
3.4	Control structures	9
3.5	Functions	9
3.6	Classes	10
3.7	Mathematical functions	10
3.8	Plotting of Functions	11
3.9	Importing data	11
4	Signal conversion	13
4.1	A/D conversion	13
4.2	D/A Conversion	13
5	Sampling of Analogue Signals	13
5.1	Normalised frequency	13
5.2	Nyquist frequency	14
5.3	Reconstruction of an analogue signal: Sampling theorem	16
6	Quantisation of analogue signals	16
6.1	Quantisation error	17
7	Frequency representation of signals	18
7.1	Continuous time and frequency	19
7.1.1	Periodic signals	19
7.1.2	A-periodic signals	20
7.2	Sampled time and/or frequency	20
7.2.1	Discrete time Fourier Transform (DTFT)	20
7.2.2	The effect of time domain sampling on the spectrum (Sampling theorem)	20
7.2.3	Discrete Fourier Transform (DFT)	23
7.2.4	Properties of the DFT	25

7.2.5	Problems with finite length DFTs	26
7.2.6	Fast Fourier Transform	26
7.3	Software	29
7.4	Visualisation	29
8	Causal Signal Processing	31
8.1	Causality	31
8.2	Convolution of Causal Signal	31
8.3	Laplace transform	32
8.4	Filters	34
8.4.1	How to characterise filters?	34
8.5	The z-transform	34
8.6	Frequency response of a sampled filter	35
8.7	FIR Filter	37
8.7.1	Generally FIR filters do not perform a convolution operation	38
8.7.2	FIR filter implementations	39
8.7.3	Fixed point FIR filters	40
8.7.4	Constant group delay or linear phase filter	41
8.7.5	Window functions	42
8.7.6	Python code: impulse response from the inverse DFT - The frequency sampling method	44
8.7.7	FIR filter design from ideal frequency response – The an- alytical way	45
8.7.8	Design steps for FIR filters	46
8.7.9	FIR filter design with Python's high level functions	47
8.7.10	Signal Detection: Matched filter	47
8.7.11	Adaptive FIR Least Mean Squares (LMS) filters	49
8.8	IIR Filter	51
8.8.1	Introduction	51
8.8.2	Determining the data-flow diagram of an IIR filter	52
8.8.3	General form of IIR filters	53
8.8.4	IIR filter topologies	53
8.8.5	Fixed point IIR filters	55
8.8.6	Filter design based on analogue filters	56
8.8.7	Adaptive IIR filter: The Kalman filter	59
8.9	The role of poles and zeros	60
8.9.1	Zeros	61
8.9.2	Poles	61
8.9.3	Stability	62
8.9.4	Design of an IIR notch filter	62
8.9.5	Identifying filters from their poles and zeroes	63

1 Acknowledgements

A big thanks to Fiona Baxter for typing up the handwritten lecture notes and turning them into \LaTeX .

2 Introduction

This handout introduces you into the basic concepts of Digital Signal Processing. In contrast to the YouTube clips it focusses more on the theoretical aspects of it and its analytical derivations. The (video-) lectures cover more practical programming examples in Python and C++.

2.1 Suggested reading

- Digital Signal Processing by John G. Proakis and Dimitris K Manolakis
- Digital Signal Processing: System Analysis and Design by Paulo S. R. Diniz, Eduardo A. B. da Silva, and Sergio L. Netto

2.2 Advantages of digital signal processing

- flexible: reconfigurable in software!
- easy storage: numbers!
- cheaper than analogue design
- very reproducible results (virtually no drift)

2.3 Development boards

Buy a DSP development board and play with it. You get them either directly from a company or from a distributor such as Farnell or RS. Also general purpose boards such as the Raspberry PI and Arduino are great platforms for DSP.

3 Python

Python is a fully featured scripting language and it's very well thought through in its approach. It has a vast library of modules which means that only rarely you need implement low level functionality. It can be programmed both in a quick

and dirty approach to try out new ideas and at the same time it supports object oriented programming so that truly professional programs can be developed.

Python has a straightforward interface to call C and C++ functions so that you can also mix it with C/C++ to speed up processing. For example a digital filter class might have the design parts written in python and the actual fast filter routines implemented in C. All this is supported by the python package management so that you can also deploy mixed python / C code without any problems. For example the AI library by google "tensorflow" can be installed as a python package and in the background it uses the GPUs on your computer without even noticing.

Here, I just comment on certain aspects of Python which are important for DSP. Check out the full documentation of python if you are interested beyond DSP in it. In particular it has also very easy to use support for GUI programming.

3.1 How to use Python?

There are two ways to use Python. The one is interactive within the python console and then the other method is by running a python script.

Development environments such as anaconda have both the console and then scripting facility integrated under one roof.

The python console The python console you can just start by typing "python":

```
Python 3.8.10 (default, Jun  2 2021, 10:49:15)
[GCC 9.4.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

then you can type in commands one by one. Under Windows you might need to point the PATH variable to your python program or you let Anaconda do it for you.

Python scripts Scripts are text files which contain python commands. One command per line and these are executed line by line.

Write your commands in any text editor and save it as a .py-file. However, it's recommended to use an editor which supports at least python's indentation which is a very important feature of python. We'll see that indentation in python defines the body of a function, the content of a loop or the code executed conditionally.

So, it's recommended for any serious work to install an editor which supports python such as spyder, VS-code, py-charm or emacs.

Once you have your script file saved just type “python myscript.py” and it will be executed.

Help There is an abundance of online help available for Python. The official Python pages have an excellent tutorial (<https://docs.python.org/3/tutorial/>) and of course you find answers to all your questions on Stack Overflow.

If you need help for a function you can use the in built help function “help(thecommand)” to get help for “thecommand”. Try out “help(print)” to get help for the “print” command.

3.2 Packages/Modules

The power of Python comes from thousands of modules which are shipped with python or can be installed via python’s package manager. A package contains usually many modules. The standard package manager of python is called “pip” which downloads the package for you and installs it. This is not just very convenient but also is safe because the download location is the official python repository.

In general you import modules with the “import” command. For example “import numpy as np” imports the module numpy into python and abbreviates it as “np” so that you’d write “np.sin(c)” to call the sine function in numpy.

The import command loads essentially a file called “numpy.py” so that you can write your own modules just by creating file with that name. See below in the “class” section for an example. There are numerous ways how to import modules but the “import as” way is by far the most popular.

Important modules are:

- **numpy**: Numerical operations for number crunching, matrix/array operations, trigonometric functions, etc. It offers fast numerical operations on arrays and matrices. It also has extensive Fourier Transform commands.
- **pylab**: Plotting functions (MATLAB style), statistics and data fitting functions
- **scipy**: Scientific computing package. In particular we are interested in `scipy.signal` which contains all the signal processing functions and filter design commands. It also offers Fourier Transforms, data io routines, vector quantisation, linear algebra, optimisation tools and statistics.

3.3 Data types

Python has all the standard data types such as int, float, string, ... and a few more. You can check which type a variable has with: `type(variable)`. Important for DSP are these types:

Numbers Numbers are represented in a standard way, including complex numbers:

- $6E23 = 6 \cdot 10^{23}$
- `6+2j` is a complex number which is just generated by adding a "j" to the imaginary part.
- `complex(0,1)` creates `1j` and is useful if you have variables which you want to combine into a complex number.

Note that python has implicit types which might make you trip over. If you write `"a=1"` then `"a"` is an integer whereas `"a=1.0"` is a floating point number. You can force a certain type by using the typecast operator, for example `"a = int(b)"` to create an integer.

Lists There are numerous data types in Python which allow storing collections of data. For DSP we need: lists, tuples and arrays. To clear up the confusion from the start we are going through them one by one. First up: lists!

- `p = [1,9,77]` is a list with the components 1, 9 and 77.
- `p[2]` returns the third element of the list (the index counts from 0). Note the index needs to be an integer so you might need to cast a floating point variables first into an integer, for example `"p[int(a)]"`.
- `p[1:5]` extracts from the list all elements from index number 1 to index number 4 (!). The index 5 is excluded! This is called splicing. Check out the online docs for more examples.
- `p = range(1,4)` generates 1, 2, 3 which is the same as a "for loop" in C: `for(int i=1; i<4; i++)`
- `p.append(10)` adds a new element to the end of the list. You see that lists are essentially classes with methods (see below)! Check out the documentation for more options to manipulate lists.
- `len(p2)` provides the number of elements.

Tupels Tupels are similar to lists but they are rather used as “containers” to package up different variables into one variable. For example, you might want to put the sampling rate and the resolution of an ADC into a tuple as they belong together.

- `p = (1,9,'kk')` is a tuple with the two integer numbers 1,2 and 'kk'.
- The round brackets are often omitted: `p = 1,9,'kk'` which implicitly defines a tuple.
- `p[1]` returns the 2nd element of the tuple (the index counts from 0).
- Extracting a tuple into separate variables also works: `a,b,c = p` assigns 1 to a, 9 to b and 'kk' to c. This is very convenient for functions if you want to return more than one result. Just package them up into a tuple and return it.

numpy Arrays/Matrices Numpy arrays or matrices are used for number crunching and are ideal for DSP. Use these for all of your DSP data manipulation. They are essentially C arrays with a thin wrapper around them. In order to create and manipulate these arrays you need to use the numpy functions.

- `a = np.array([1, 2, 5])` which creates a numpy array from a standard Python list containing 1,2,5.
- `b = np.array([[6, 2], [9, 7]])` creates a two dimensional array.
- There are a lot of convenience functions in numpy to create arrays just containing zeros, ones or ascending numbers etc. Check out the numpy docs.

Since numpy arrays are close relatives to C arrays they have an additional field which is called “dtype” which determines the type of the array elements, for example “int16” (a 16 bit wide integer).

```
>>> data
array([-12288, -12545, -12798, ...,   511,   513,   513], dtype=int16)
>>> data.dtype
dtype('int16')
>>> data.astype(float)
array([-12288., -12545., -12798., ...,   511.,   513.,   513.] )
```

where the member function “astype” can convert from one array type to another, here from int16 to floating point.

3.4 Control structures

Conditionals Conditional operation is done in exactly the same way as in other languages:

```
if a > c:
    a = c
    c = c*2
print(c)
```

The indent of 4 characters indicates the conditionally executed commands. The same applies for do and while loops where the body has an indent. Here, `print(c)` is always executed because it doesn't have an indent.

Loops Loops in python are always done with the help of iterators. This means that we progress through the elements of a list. Imagine you want to calculate the average of the array `y`.

```
avg = 0
for i in range(len(y)):
    avg = avg + y[i]
    print("momentary average=", avg)
print(avg/len(y))
```

The “for” statement iterates through all elements of an array. In this case we have created a list with the “range” command. Important: the *indent* indicates the commands which the loop iterates through. This is the general rule in python. There are also while/do loops available.

3.5 Functions

Function definitions can happen anywhere in the program or in a separate file. For example, the above average calculations can be defined as:

```
def myaverage(g):
    avg = 0
    for i in range(len(y)):
        avg = avg + y[i]
    return (avg/len(y))
```

and then called from the main program as `myaverage(y)`. The keyword “def” indicates the start of the function definition. Note again the 4 character indent for the function body and the colon after the function name. This is characteristic for python.

3.6 Classes

Classes contain functions and variables. They usually have a constructor called `init` which is a special function which is called when an instance of a class is created.

```
class MyAmazingClass:
    def __init__(self,coeff):
        # Your init code here.
        self.mylocalcoeff = coeff

    def doSomeProcessing(self,data):
        # we process data
        data = data * self.mylocalcoeff
        return data
```

(!) Note the special variable “self” which is *the first* argument of every method in a class. Variables which are part of the class are references by “self”. For example the constructor argument “coeff” is saved in the class as “mylocalcoeff” and is then available throughout the lifetime of the class as “self.mylocalcoeff”. When we call the method “doSomeProcessing” then “self.mylocalcoeff” is multiplied with the “data” argument and the result is returned. Anything which hasn’t got “self.” in front of it will be forgotten after the function has terminated. In this way you distinguish which variables are kept for the lifetime of the class or are just used temporarily.

Classes are just types and need to be instantiated:

```
f = MyAmazingClass(5)
a = f.doSomeProcessing(10)
b = f.doSomeProcessing(20)
```

which will give us 50 as a result for the variable “a” and 100 for the variable “b”. The variable “f” is the instance of the class `MyAmazingClass` and the argument initialises `self.mylocalcoeff` with 5.

3.7 Mathematical functions

Just import the `numpy` library which contains all the mathematical functions and constants such as `pi`. `import numpy as np`. It has all the standard functions such as `np.sin`, `np.cos`, `np.tan`, `np.exp`, ...

Note that these functions also take arrays as arguments. For example, inputting the array `x` into the sine:

```
x = np.array(range(6))
y = np.sin(x);
```

yields again an array in *y*! All elements have been processed at the *same* time.

3.8 Plotting of Functions

PyLab has a vast amount of plotting functions which supports simple plotting with just two commands (plot and show) up to complex animations of 3D data. Plotting a sine wave can be achieved with:

```
import matplotlib.pyplot as plt
import numpy as np
x = np.linspace(0,2*np.pi,100)
y = np.sin(x)
plt.plot(x,y)
plt.show()
```

where I have used the numpy linspace command to create an array which runs from 0 to two pi with 100 samples.

If you want to create more plot-windows you can use the command `plt.figure`. Access the windows with the commands `plt.figure(1)`, `plt.figure(2)`.

It is also possible to combine plot windows into one with the command `subplot`.

How to make it look nice? Here is an example:

```
plt.xlabel('time (s)')
plt.ylabel('voltage (V)')
plt.title('That is an AC voltage over time!')
plt.grid(True)
```

Generally the best approach is to use an example and then customise it to your needs.

3.9 Importing data

Space/Tab/Comma separated data Numpy has a very handy import function for data as long as the the elements in the data file are separated with spaces, tabs or commas. To load such data files just type in: `x = np.loadtxt('mydatafile.dat')`. This creates a two dimensional array which reflects exactly the structure in the datafile.

Thus, if you want to write software which creates Python readable data then just export the data as space/tab separated ASCII files.

WAV files Scipy contains in its submodule io a subsubmodule wavfile which supports reading and writing wav files:

```
import scipy.io.wavfile as wavfile
fs,data = wavfile.read('kongas.wav');
```

fs,data is a tuple containing the sampling rate and then the data.

```
import scipy.io.wavfile as wavfile
wavfile.write('kongas_filt.wav',fs,y3);
```

whereas 'fs' is the sampling rate and y3 the data – usually normalised between -1 and +1 for floating point numbers.

4 Signal conversion

$$\text{Analogue} \rightarrow \text{A/D} \rightarrow \text{Digital processing} \rightarrow \text{D/A} \rightarrow \text{Analogue} \quad (1)$$

4.1 A/D conversion

1. sampling

$$\underbrace{X_a(nT)}_{\text{analogue signal}} \equiv \underbrace{x(n)}_{\text{discrete data}} \quad (2)$$

- X_a : analogue signal
- T : sampling interval, $\frac{1}{T} = F_s$ is the sampling rate
- $x(n)$: discrete data

2. **quantisation**: Continuous value into discrete steps. $X(n) \rightarrow X_q(n)$

3. **Coding** $X_q(n) \rightarrow$ into binary sequence or integer numbers.

4.2 D/A Conversion

Back to analogue. The D/A converter needs to interpolate between the samples to get the original signal back without error! This is usually done by a simple low pass filter.

5 Sampling of Analogue Signals

The analogue signal X_a is sampled at intervals T to get the sampled signal $X(n)$.

$$X(n) = X_a(nT), \quad -\infty < n < +\infty \quad (3)$$

5.1 Normalised frequency

Note that the sampling frequency

$$F_s = \frac{1}{T_s} \quad (4)$$

is lost and needs to be stored additionally so that the analogue signal can be reconstructed:

$$X_a(t) = X_a(nT) = X_a\left(\frac{n}{F_s}\right) \quad (5)$$

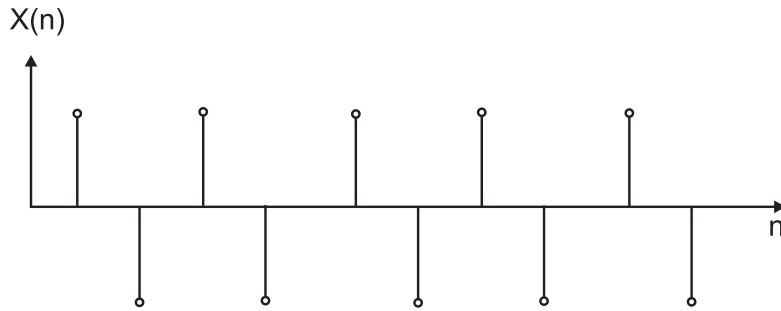


Figure 1: Maximum frequency of a sampled signal is half the normalised frequency.

What is the frequency range of our digital signal? What is the maximum frequency in digital terms? Fig 1 shows that the max frequency is 0.5 because we need two samples to represent a “wave”. The frequency range $0 \dots 0.5$ is called the *normalised frequency* range. What is the relation between normalised frequency and sampling rate?

Let's have a look at the analog signal $X_a(t)$ with the frequency F :

$$X_a(t) = A \cos(2\pi Ft) \quad (6)$$

and we sample it only at:

$$t = n/F_s \quad (7)$$

Then the digital signal is:

$$X_a(n) = A \cos(2\pi F/F_s n) \quad (8)$$

Now, we can define the normalised frequency as:

$$\text{Normalised frequency: } f = \frac{F}{F_s} \quad (9)$$

which has its max value at 0.5 which represents one period of a sine wave within two samples¹

5.2 Nyquist frequency

Recall that the max value for the normalised frequency f is 0.5 and that:

$$X_a(n) = A \cos(2\pi n f) \quad (10)$$

¹Note that the normalised frequency in Python's scipy is differently defined. In scipy it is annoyingly defined as $f_{\text{scipy}} = 2\frac{F}{F_s}$. So, in other words $f_{\text{scipy}} = 1$ is the Nyquist frequency instead of $f = 0.5$ as normally defined.

with n as an integer because we are sampling.

What happens above $f > 0.5$? Imagine $f = 1$

$$X_a(n) = \cos(2\pi n) = 1 \quad f = 1 \quad (11)$$

which gives us DC or zero frequency. The same is true for $f = 2, 3, 4, \dots$. We see that above $f = 0.5$ we never get higher frequencies. Instead they will always stay between $0 \dots 0.5$ for the simple reason that it is not possible to represent higher frequencies. This will be discussed later in greater detail.

The ratio F/F_s must be lower than 0.5 to avoid ambiguity or in other words the maximum frequency in a signal must be lower than $\frac{1}{2}F_s$. This is the Nyquist frequency.

If there are higher frequencies in the signal then these frequencies are “folded down” into the frequency range of $0 \dots \frac{1}{2}F_s$ and creating an alias of its original frequency in the so called “baseband” ($f = 0 \dots 0.5$). As long as the alias is not overlapping with other signal components in the baseband this can be used to downmix a signal. This leads to the general definition of the sampling theorem which states that the bandwidth B of the input signal must be half of the sampling rate F_s :

$$B < \frac{1}{2}F_s \quad (12)$$

The frequency $\frac{1}{2}F_s$ is called the Nyquist frequency.

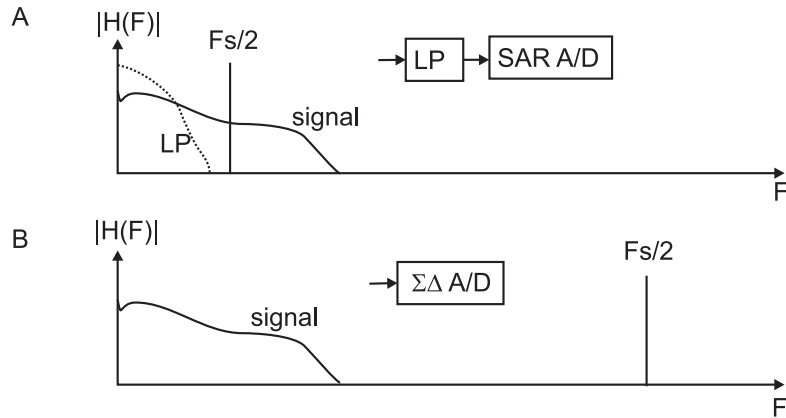


Figure 2: Anti alias filtering. A) with a lowpass filter. B) with a sigma delta converter with very high sampling rate.

What do we do if the signal contains frequencies above $\frac{1}{2}F_s$? There are two ways to tackle this problem: The classical way is to use a lowpass filter (see Fig. 2A) which filters out all frequencies above the Nyquist frequency. However

this might be difficult in applications with high resolution A/D converters. Alternatively one can use a much higher sampling rate to avoid aliasing. This is the idea of the sigma delta converter which operates at sampling rates hundred times higher than the Nyquist frequency.

5.3 Reconstruction of an analogue signal: Sampling theorem

Is it possible to reconstruct an analogue signal from a digital signal which contains only frequencies below the Nyquist frequency?

$$F_s > 2F_{\max} \quad (13)$$

where F_{\max} is max frequency in the signal which we represent by sine waves:

$$x_a(t) = \sum_{i=1}^n A_i \cos(2\pi F_i t + \Theta_i) \quad (14)$$

The analogue signal $x_a(t)$ can be completely reconstructed if:

$$g(t) = \frac{\sin 2\pi B t}{2\pi B t} \quad (15)$$

with

$$B = F_{\max} \quad (16)$$

$$x_a(t) = \sum_{h=-a}^a x_a\left(\frac{n}{F_s}\right) g\left(t - \frac{h}{F_s}\right) \quad (17)$$

6 Quantisation of analogue signals

A/D converters have a certain resolution. For example, the MAX1271 has a resolution of 12 bits which means that it divides the input range into 4096 equal steps (see Fig 3).

$$\Delta = \text{quantisation step} = \frac{x_{\max} - x_{\min}}{L - 1} \quad (18)$$

where $x_{\max} - x_{\min}$ is the dynamic range in volt (for example 4.096V) and L is the number of quantisation steps (for example, 4096). Δ is the quantisation step which defines minimal voltage change which is needed to see a change in the output of the quantiser. The operation of the quantiser can be written down as:

$$x_q(n) = Q[x(n)] \quad (19)$$

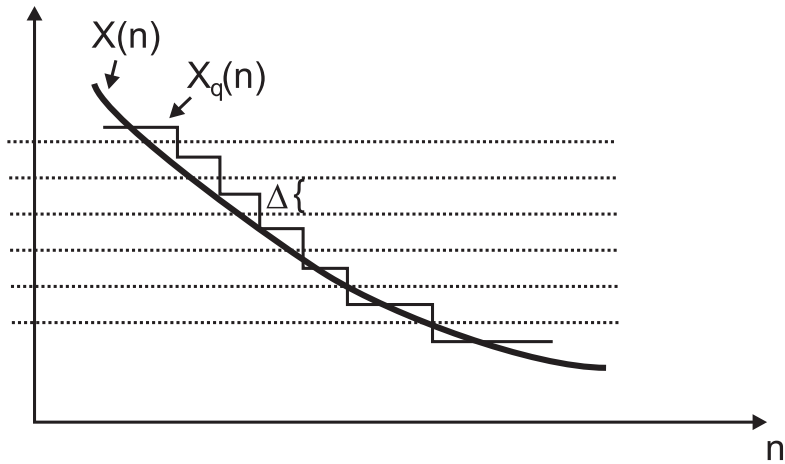


Figure 3: Δ is the quantisation step.

6.1 Quantisation error

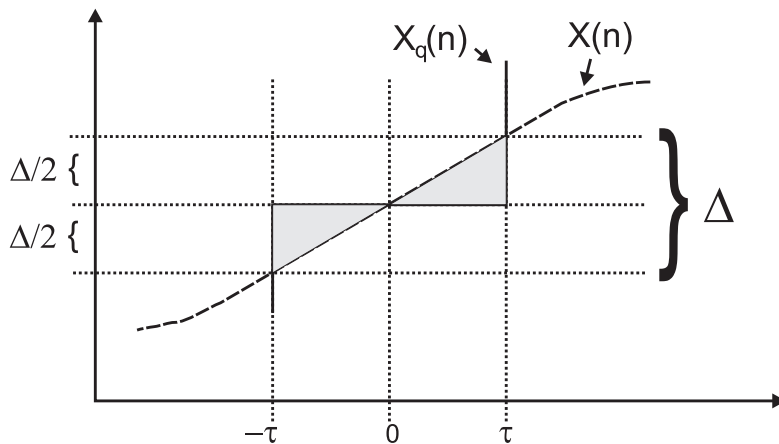


Figure 4: Illustration of the quantisation error. It is zero at $t = 0$ and increases to the edges of the sampling interval. Illustrated is the worst case scenario. This repeats in the next sampling interval and so forth.

Fig. 4 shows error produced by the quantiser in the worst case scenario. From that it can be seen that the maximum quantisation error is half the quantisation step:

$$-\frac{\Delta}{2} \leq e(n) \leq \frac{\Delta}{2} \quad (20)$$

The smaller the quantisation step Δ the lower the error!

What is the mean square error P_q ?

$$P_q = \frac{1}{\tau} \int_0^\tau e_q^2(t) dt \quad (21)$$

$$P_q = \frac{1}{\tau} \int_0^\tau \left(\frac{\Delta}{2\tau} \right)^2 t^2 dt \quad (22)$$

$$= \frac{\Delta^2}{4\tau^3} \int_0^\tau t^2 dt \quad (23)$$

$$P_q = \frac{\Delta^2}{12\tau^3} \tau^3 = \frac{\Delta^2}{12} \quad (24)$$

What is the relative error to a sine wave?

$$P_x = \frac{1}{T_p} \int_0^{T_p} (A \cos \Omega t)^2 dt = \frac{A^2}{2} \quad (25)$$

Ratio to signal power to noise:

$$\text{SQNR} = \frac{P_x}{P_q} = \frac{A^2}{2} \cdot \frac{12}{\Delta^2} \quad (26)$$

$$= \frac{6A^2}{\Delta^2} \quad (27)$$

This equation needs to be interpreted with care because increasing the amplitude of the input signal might lead to saturation if the input range of the A/D converter is exceeded.

7 Frequency representation of signals

Often we are more interested in the frequency representation of signals than the evolution in time.

We will use small letters for time domain representation and capital letters for the frequency representation, for example $X(F)$ and $x(t)$.

7.1 Continuous time and frequency

7.1.1 Periodic signals

Periodic signals can be composed of sine waves :

$$x(t) = \sum_{k=-\infty}^{\infty} c_k e^{j2\pi k F_1 t} \quad (28)$$

where F_1 is the principle or fundamental frequency and $k \neq 1$ are the harmonics with strength c_k . Usually c_1 is set to 1. This is a Fourier series with c_k as coefficients. For example an ECG has a fundamental frequency of about 1Hz (60 beats per minute). However, the harmonics give the ECG its characteristic peak like shape.

How do we find the coefficients c_k ?

$$c_k = \frac{1}{T_p} \int_{T_p} x(t) e^{-j2\pi k F_1 t} dt \quad (29)$$

For simple cases there are analytical solutions for c_k , for example for square waves, triangle wave, etc.

What are the properties of c_k ?

$$c_k = c_{-k}^* \Leftrightarrow x(t) \text{ is real} \quad (30)$$

or

$$c_k = |c_k| e^{j\theta_k} \quad (31)$$

$$c_{-k} = |c_k| e^{-j\theta_k} \quad (32)$$

Proof: with the handy equation...

$$\cos z = \frac{1}{2} (e^{zi} + e^{-zi}) \quad (33)$$

we get

$$x(t) = c_0 + \sum_{k=1}^{\infty} |c_k| e^{j\theta_k} e^{j2\pi k F_1 t} + \sum_{k=1}^{\infty} |c_k| e^{-j\theta_k} e^{-j2\pi k F_1 t} \quad (34)$$

$$= c_0 + 2 \sum_{k=1}^{\infty} |c_k| \cos(2\pi k F_1 t + \theta_k) \quad (35)$$

How are the frequencies distributed? Let's have a look at the frequency spectrum of a periodic signal: $P_k = |c_k|^2$

- There are discrete frequency peaks
- Spacing of the peaks is $\frac{1}{T_1} = F_1$
- Only the positive frequencies are needed: $c_{-k} = c_k^*$

7.1.2 A-periodic signals

In case nothing is known about $X(t)$ we need to integrate over all frequencies instead of just the discrete frequencies.

$$X(F) = \int_{-\infty}^{\infty} x(t)e^{-j2\pi Ft} dt \quad (36)$$

Consequently, the frequency spectrum $X(F)$ is continuous.

$$x(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} X(F)e^{j2\pi Ft} dF \quad (37)$$

7.2 Sampled time and/or frequency

7.2.1 Discrete time Fourier Transform (DTFT)

The signal $x(n)$ is discrete whereas the resulting frequency spectrum is considered as continuous (arbitrary signals).

- Analysis or direct transform:

$$X(\omega) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\omega n} \quad (38)$$

- Synthesis or inverse transform:

$$x(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} X(\omega)e^{j\omega n} d\omega \quad (39)$$

$$= \frac{1}{2\pi} \int_{-0.5}^{0.5} X(f)e^{j2\pi fn} dt \quad (40)$$

note the range here. f is the normalised frequency.

7.2.2 The effect of time domain sampling on the spectrum (Sampling theorem)

What effect has time domain sampling on the frequency spectrum²? Imagine we have an analog spectrum $X_a(F)$ from a continuous signal $x(t)$. We want to know how the spectrum $X(F)$ of the discrete signal $x(n)$ looks like.

$$X(F) \Leftrightarrow X_a(F) \quad (41)$$

²This derivation is loosely based on [Proakis and Manolakis \(1996\)](#)



Figure 5: Effect of sampling on the spectrum.

The signal is represented by $x(n)$ so that we can equate the Fourier transforms of the sampled spectrum $X(F)$ and of the analogue spectrum $X_a(F)$.

$$\int_{-0.5}^{0.5} \underbrace{X(f)}_{\text{sampled}} e^{j2\pi f n} df = \int_{-\infty}^{+\infty} \underbrace{X_a(F)}_{\text{cont}} e^{j2\pi n F / F_s} dF \quad (42)$$

Obviously $X(f)$ must be different to accommodate the different integration ranges. The trick is now to divide the integral on the right hand side of Eq. 42 into chunks to make it compatible to the range on the left hand side.

Remember that the normalised frequency is $f = F/F_s$ which allows us to change the integration to analogue frequency on both sides:

$$\frac{1}{F_s} \int_{-F_s/2}^{F_s/2} X\left(\frac{F}{F_s}\right) e^{j2\pi n F / F_s} dF = \int_{-\infty}^{+\infty} X_a(F) e^{j2\pi n F / F_s} dF \quad (43)$$

and now we divide the right hand side into chunks of F_s which corresponds to the integration range on the left hand side.

$$\int_{-\infty}^{+\infty} X_a(F) e^{j2\pi n \frac{F}{F_s}} dF = \sum_{k=-\infty}^{\infty} \int_{-\frac{1}{2}F_s + kF_s}^{+\frac{1}{2}F_s + kF_s} X_a(F) e^{j2\pi n \frac{F}{F_s}} dF \quad (44)$$

$$= \sum_{k=-\infty}^{\infty} \int_{-\frac{1}{2}F_s}^{+\frac{1}{2}F_s} X_a(F - kF_s) \underbrace{e^{j2\pi n \frac{F}{F_s}}}_{kF_s \text{ omit.}} dF \quad (45)$$

$$= \int_{-\frac{1}{2}F_s}^{+\frac{1}{2}F_s} \underbrace{\sum_{k=-\infty}^{\infty} X_a(F - kF_s)}_{=X(F) \text{ of Eq. 43}} e^{j2\pi n \frac{F}{F_s}} dF \quad (46)$$

This gives us now an equation for the sampled spectrum:

$$X(F/F_s) = F_s \sum_{k=-\infty}^{\infty} X_a(F - kF_s) \quad (47)$$

$$X(f) = F_s \sum_{k=-\infty}^{\infty} X_a[(f - k)F_s] \quad (48)$$

This equation can now be interpreted. In order to get the sampled spectrum $X(F)$ we need to make copies of the analog spectrum $X_a(F)$ and place these copies at multiples of the sampling rate F_s (see Fig. 5). This illustrates also the *sampling theorem*: if the bandwidth of the spectrum is wider than $F_s/2$ then the copies of the analogue spectrum will overlap and reconstruction would be impossible. This is called aliasing. Note that it is not necessary bad that the spectrum of the analogue signal lies within the range of the so called “base band” $-F/2 \dots F/2$. It can also lie in another frequency range further up, for example $-F/2 + 34 \dots F/2 + 34$ as long as the *bandwidth* does not exceed $F_s/2$. If it is placed further up it will automatically show up in the baseband $-F/2 \dots F/2$ which is called “fold down”. This can be used for our purposes if we want to down mix a signal.

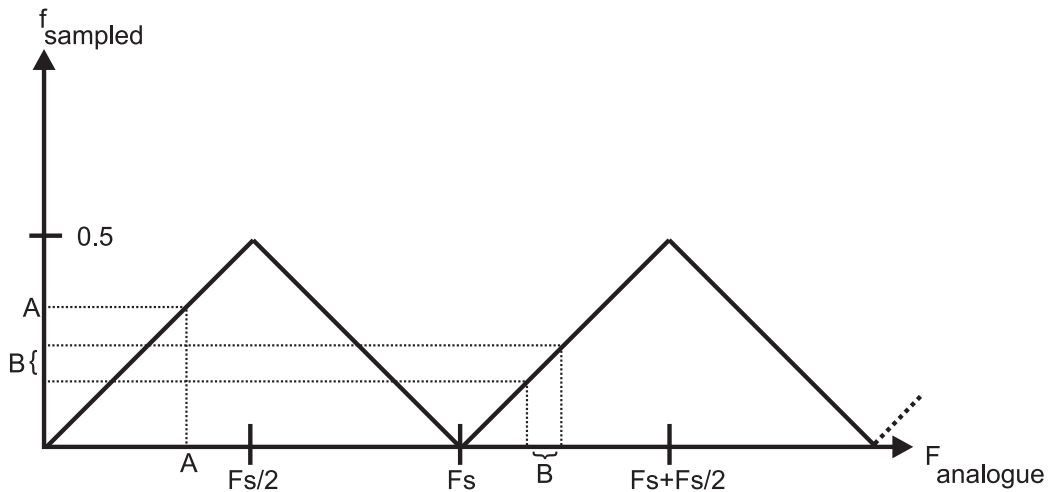


Figure 6: Mapping of frequencies in the analogue domain (F) and sampled domain (f). F_s is the sampling rate.

With the insight from these equations we can create a plot of how analogue frequencies map onto sampled frequencies. Fig 6 shows how the analogue frequencies F_{analogue} map on the normalised frequencies f_{sampled} . As long as the analogue frequencies are below $F_s/2$ the mapping is as usual as shown in Fig 6A. Between $F_s/2$ and F_s we have an inverse mapping: an increase in analogue frequency causes a decrease in frequencies. Then, from F_s we have again an increase in frequencies starting from DC. So, in general if we keep a bandlimited signal within one of these slopes (for example from $F_s \dots F_s + 1/2F_s$ as shown in Fig 6B) then we can reproduce the signal.

This leads us to the generalised Nyquist theorem: if a bandpass filtered signal has a bandwidth of B then the minimum sampling frequency is $F_s = 2B$.

7.2.3 Discrete Fourier Transform (DFT)

So far we have only sampled in the time domain. However, on a digital computer the Fourier spectrum will always be a discrete spectrum.

The discrete Fourier Transform (DFT) is defined as:

$$X(k) = \sum_{n=0}^{N-1} x(n)e^{-j2\pi kn/N} \quad k = 0, 1, 2, \dots, N-1 \quad (49)$$

where N is the number of samples in both the time and frequency domain.

The inverse discrete Fourier Transform (IDFT) is defined as:

$$x(n) = \frac{1}{N} \sum_{k=0}^{N-1} X(k)e^{j2\pi kn/N} \quad n = 0, 1, 2, \dots, N-1 \quad (50)$$



Figure 7: Effect of sampling in the frequency domain. The inverse in the time domain $x(n)$ contains repetitions of the original signal.

What is the effect in the timedomain of this discretisation³? We start with the continuous Fourier transform and discretise it into N samples in the frequency

³This derivation is loosely based on [Proakis and Manolakis \(1996\)](#)

domain:

$$X\left(\frac{2\pi}{N}k\right) = \sum_{n=-\infty}^{\infty} x(n)e^{-j\frac{2\pi}{N}kn} \quad k = 0, \dots, N-1 \quad (51)$$

Let's subdivide the sum into chunks of length N :

$$X\left(\frac{2\pi}{N}k\right) = \sum_{l=-\infty}^{\infty} \sum_{n=lN}^{lN+N-1} x(n)e^{-j\frac{2\pi}{N}kn} \quad (52)$$

$$= \sum_{l=-\infty}^{\infty} \sum_{n=0}^{N-1} x(n-lN)e^{-j\frac{2\pi}{N}kn} \quad (53)$$

$$= \sum_{n=0}^{N-1} \underbrace{\sum_{l=-\infty}^{\infty} x(n-lN)}_{\text{Periodic repetition!}} e^{-j\frac{2\pi}{N}kn} \quad (54)$$

We note the following:

- Ambiguity in the time domain
- The signal is repeated every N samples

Practically this repetition won't show up because the number of samples is limited to N in the inverse transform. However, for operations which shift signals in the frequency domain it is important to remember that we shift a periodic time series. If we shift it out at the end of the array we will get it back at the start of the array.

Fig. 8 shows an example of a DFT. It is important to know that the spectrum is mirrored around $N/2$. DC is represented by $X(0)$ and the Nyquist frequency $F_s/2$ is represented by $X(N/2)$. The mirroring occurs because the input signal $x(n)$ is *real*. This is important if one wants to modify the spectrum $X(F)$ by hand, for example to eliminate 50Hz noise. One needs to zero two elements of $X(F)$ to zero. This is illustrated in this python code:

```
import scipy as sp
yf=sp.fft(y)
# the sampling rate is 1kHz. We've got 2000 samples.
# midpoint at the ifft is 1000 which corresponds to 500Hz
# So, 100 corresponds to 50Hz
yf[99:101+1]=0;
# and the mirror
yf[1899:1901+1]=0;
#
yi=sp.ifft(yf);
```

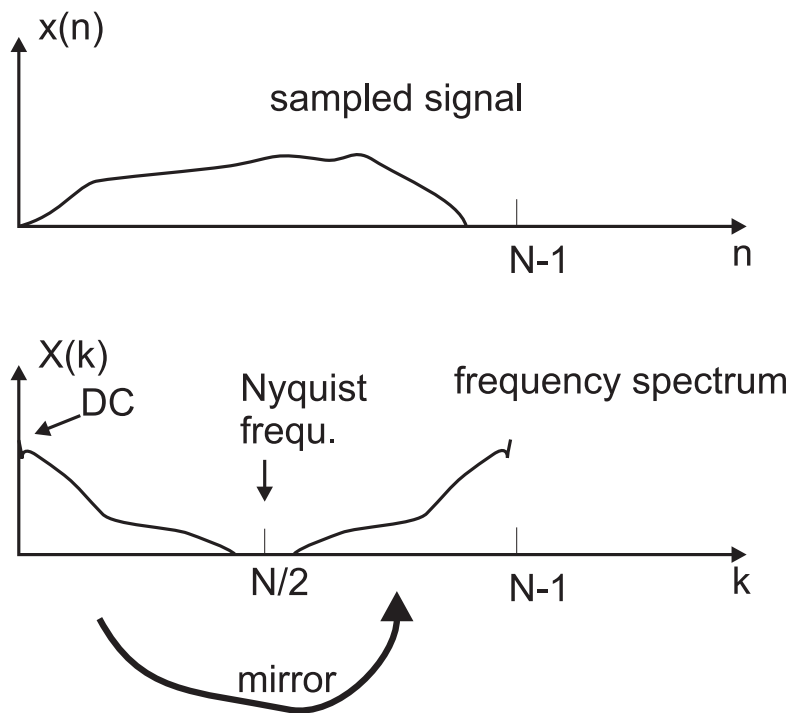



Figure 8: Example of a DFT. The sampled signal $x(n)$ is converted to the spectrum $X(k)$. Both have N samples. Note the redundancy: the spectrum is mirrored around $N/2$. The first value $X(0)$ is the DC value of the signal. However this is the only sample which is not mirrored.

This filters out the 50Hz hum from the signal y with sampling rate 1000Hz. The signal y_i should be real valued again or contain only very small complex numbers due to numerical errors.

7.2.4 Properties of the DFT

- Periodicity:

$$x(n + N) = x(n) \quad (55)$$

$$X(k + N) = X(k) \quad (56)$$

- Symmetry: if $x(n)$ real:

$$x(n) \text{ is real} \Leftrightarrow X^*(k) = X(-k) = X(N - k) \quad (57)$$

This is important when manipulating $X(k)$ by hand.

- Time Reversal:

$$x(n) \leftrightarrow X(k) \quad (58)$$

$$x(-n) \leftrightarrow X(N - k) \quad (59)$$

- Circular convolution:

$$X_1(k)X_2(k) \leftrightarrow x_1(n) * x_2(n) \quad (60)$$

with

$$x_3(m) = \sum_{n=0}^{N-1} x_1(n)x_2(m - n) \quad (61)$$

More useful equations are at [Proakis and Manolakis \(1996, pp.415\)](#).

7.2.5 Problems with finite length DFTs

$$x_w(n) = x(n) \cdot w(n) \quad (62)$$

where x is the input signal and $w(n)$ represents a function which is 1 from $n = 0$ to $n = L - 1$ so that we have L samples from the signal $x(n)$.

To illustrate the effect of this finite sequence we introduce a sine wave $x(n) = \cos \omega_0 n$ which has just two peaks in a proper Fourier spectrum at $-\omega$ and $+\omega$.

The spectrum of the rectangular window with the width L is:

$$W(\omega) = \frac{\sin(\omega L/2)}{\sin(\omega/2)} e^{-j\omega(L-1)/2} \quad (63)$$

The resulting spectrum is then:

$$X_3(\omega) = X(\omega) * W(\omega) \quad (64)$$

Because the spectrum of $X(\omega)$ consists of just two delta functions the spectrum $X_3(\omega)$ contains the window spectrum $W(\omega)$ twice at $-\omega$ and $+\omega$ (see Fig. 9). This is called *leakage*. Solutions to solve the leakage problem? Use Windows with a narrower spectrum and with less ripples (see FIR filters).

7.2.6 Fast Fourier Transform

We can rewrite the DFT (Eq. 49) in a slightly more compact form:

$$X(k) = \sum_{n=0}^{N-1} x(n)W_N^{kn} \quad (65)$$

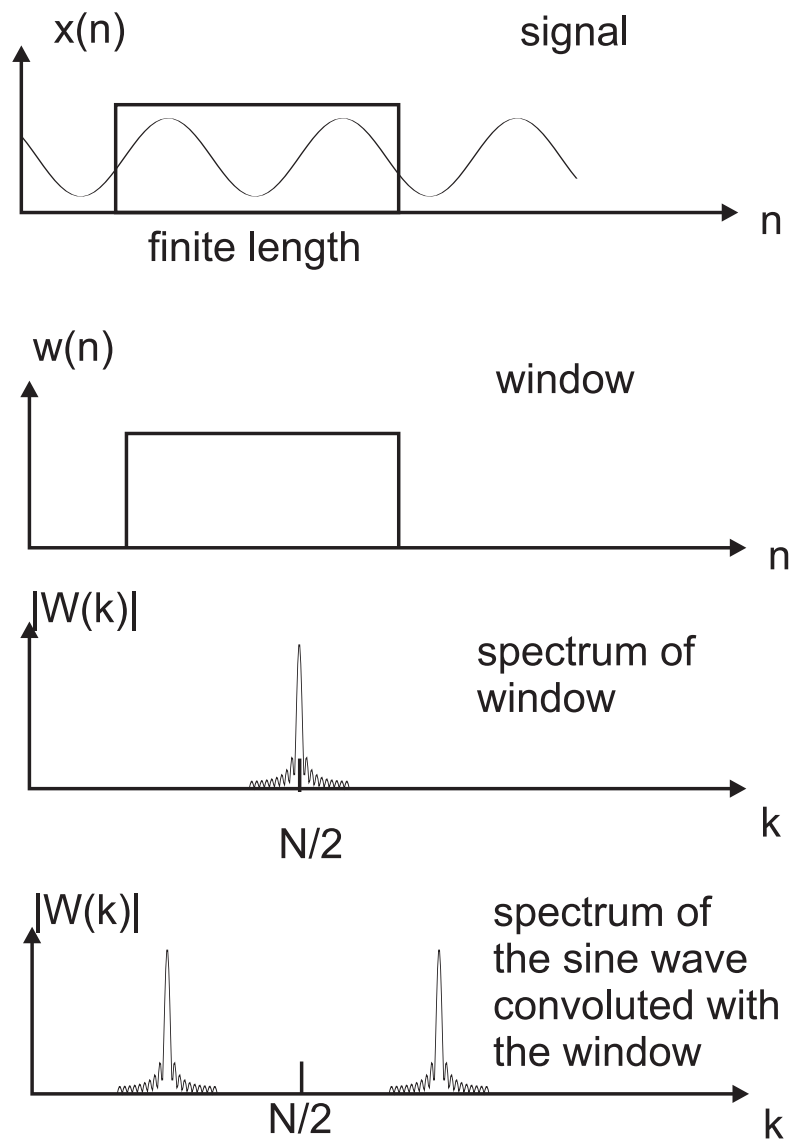


Figure 9: The effect of windowing on the DFT.

with the constant:

$$W_N = e^{-j2\pi/N} \quad (66)$$

The problem with the DFT is that it needs N^2 multiplications. How can we reduce the number of multiplications? Idea: Let's divide the DFT in an odd and an even sequence:

$$x(2m) \quad (67)$$

$$x(2m+1), \quad m = 0, \dots, \frac{N}{2} - 1 \quad (68)$$

which gives us with the trick $W_N^{2mk} = W_{N/2}^{mk}$ because of the definition Eq. 66.

$$X(k) = \sum_{m=0}^{N/2-1} x(2m)W_N^{2mk} + \sum_{m=0}^{N/2-1} x(2m+1)W_N^{k(2m+1)} \quad (69)$$

$$= \sum_{m=0}^{N/2-1} x(2m)W_{N/2}^{mk} + W_N^k \sum_{m=0}^{N/2-1} x(2m+1)W_{N/2}^{mk} \quad (70)$$

$$= F_e(k) + W_N^k F_o(k) \quad (71)$$

F_e and F_o have both half the length of the original sequence and need only $(N/2)^2$ multiplication, so in total $2 \cdot (N/2)^2 = \frac{N^2}{2}$. Basically by dividing the sequence in even and odd parts we can reduce the number of multiplications by 2. Obviously, the next step is then to subdivide the two sequences $F_e(k)$ and $F_o(k)$ even further into something like $F_{ee}(k)$, $F_{eo}(k)$, $F_{oe}(k)$ and $F_{oo}(k)$.



Figure 10: Illustration of the Fast Fourier Algorithm. The sequence of N samples is recursively divided into subsequences of odd and even samples.

In general the recipe for the calculation of the FFT is:

$$X_i(k) = X_{ie}(k) + W_L^k X_{io}(k) \quad (72)$$

W_L^k is the phase factor in front of the odd sequence. This is continued until we have only two point ($N = 2$) DFTs (see Eq. 65):

$$\text{DC:} \quad X(0) = x(0) + \underbrace{W_2^0}_1 x(1) = x(0) + x(1) \quad (73)$$

$$\text{Nyquist frequ.:} \quad X(1) = x(0) + \underbrace{W_2^1}_{-1} x(1) = x(0) - x(1) \quad (74)$$

A two point DFT operates only on two samples which can represent only two frequencies: DC and the Nyquist frequency which makes the calculation trivial. Eq. 73 is an averager and Eq. 74 is basically a differentiator which gives the max output for the sequence $1, -1, 1, -1, \dots$. Fig. 10 illustrates how to divide the initial sequence to arrive at 2 point DFTs. In other words the calculation of the full DFT is done by first calculating $N/2$ 2 point DFTs and recombining the results with the help of Eq. 72. This is sometimes called the “Butterfly” algorithm because the data flow diagram can be drawn as a butterfly. The number of complex multiplications reduces in this approach to $N \log_2 N$ which is actually the worst case scenario because many W_L^k usually turn out to be $1, -1, j, -j$ which are just sign inversions or swaps of real and imaginary parts. A clever implementation of this algorithm will be even faster.

In summary the idea behind the FFT algorithms is to divide the sequence into subsequences. Here we have presented the most popular radix 2 algorithm. The radix 4 is even more efficient and there are also algorithms for divisions into prime numbers and other rather exotic divisions. However, the main idea is always the same: subsample the data in a clever way so that the final DFT becomes trivial.

7.3 Software

In Teukolsky et al. (2007) you’ll find highly optimised C code for Fourier transforms. Most Linux distros (Ubuntu, Suse, RedHat, ...) come with an excellent FFT library called `libfftw3`.

7.4 Visualisation

The Fourier spectrum is a complex function of frequency which cannot be displayed in a convenient way. Usually the amplitude or magnitude of the spectrum is of interest (see Fig. 11) which is calculated as the absolute value $|X(F)|$ of the Fourier spectrum $X(F)$. In addition the phase might be of interest which is calculated as $\arg(X(F))$.

- Magnitude or Amplitude: $|X(F)|$
- Phase: $\arg(X(F))$



Figure 11: An ECG signal in the time- and frequency domain.

8 Causal Signal Processing

In realtime systems and any systems processing data as it arrives the future samples are not known. This is as in an analogue system where we send for example a signal into an amplifier and it needs to process it as it arrives.

However, the Fourier transform is *not real time*. We need the whole signal from the first to the last sample. For that reason we need to develop a new mathematical framework which we call “Causal Signal Processing”. These systems should ideally react as an analogue system, namely as fast as possible with little latency as possible.

In the next section we are now developing a mathematical framework for causal digital signal processing. We draw here heavily from analogue circuit design as this is by its own nature a realtime causal processor.

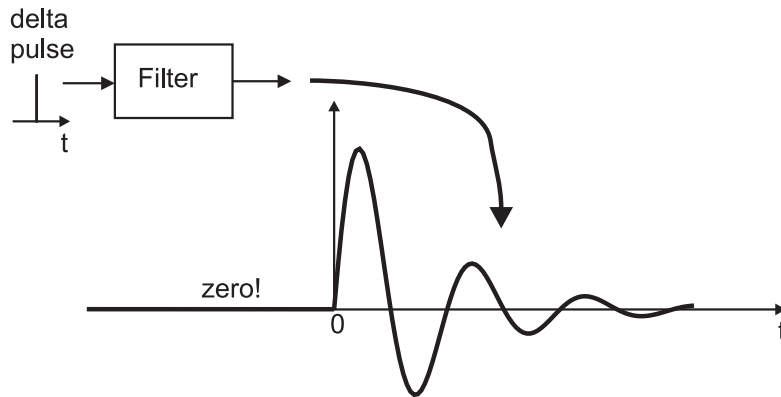


Figure 12: A causal system only *reacts* to it's input. Causal signals only evolve in positive time. Per definition the signal is zero for $t < 0$.

8.1 Causality

Fig. 12 illustrates the concept of causality. Causal systems cannot look into the future. They can only react to a certain input. Causal signals are kept zero for $t < 0$ per definition.

8.2 Convolution of Causal Signal

After having defined causality we can define the convolution:

$$y(t) = h(t) * x(t) = \int_{-\infty}^{\infty} h(t - \tau)x(\tau)d\tau \quad (75)$$

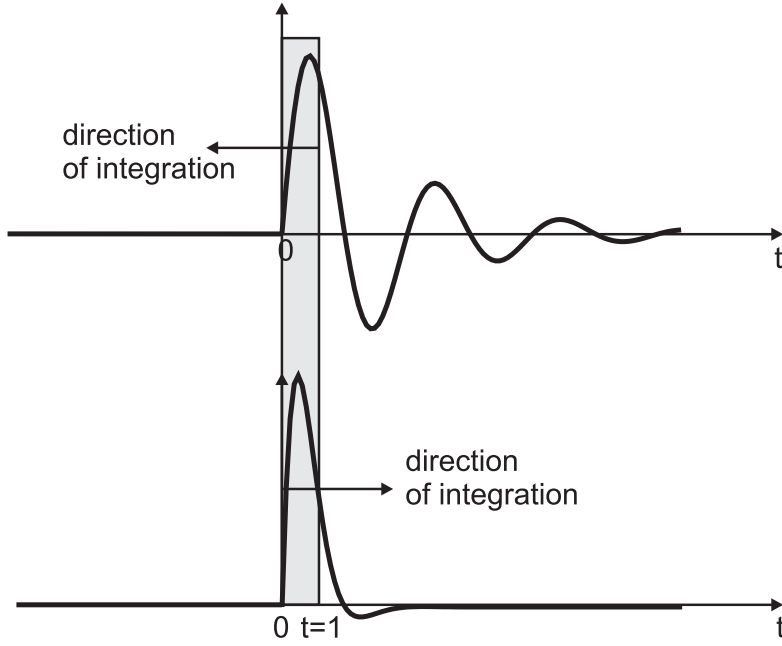


Figure 13: Illustration of the convolution. The shaded area shows the integration for $t = 1$.

$$y(n) = h(n) * x(n) = \sum_{n=-\infty}^{\infty} h(n)x(m-n) \quad (76)$$

Note the reversal of integration of the function h . This is characteristic of the convolution. At time $t = 0$ only the values at $h(0)$ and $x(0)$ are evaluated (see Fig. 13). Note that both functions are zero for $t < 0$ (causality!). At $t > 0$ the surface which is integrated grows as shown in Fig. 13 for $t = 1$.

What happens if $x(t) = \delta(t)$? Then Eq. 76:

$$y(t) = \int_{-\infty}^{\infty} h(t-\tau)\delta(\tau)d\tau = h(t) \quad (77)$$

provides us with the function $h(t)$ itself. This will be used later to determine the impulse response of the filter.

8.3 Laplace transform

The Fourier transform is not suitable for causal systems because it requires the whole signal from $-\infty < t < +\infty$. What we need is a transform which works with causal signals. This is the Laplace transform:

$$\mathbf{LT}(h(t)) = H(s) = \int_0^{\infty} h(t)e^{-st}dt \quad (78)$$

The Laplace transform has a couple of very useful properties:

- Integration:

$$\int f(\tau) d\tau \Leftrightarrow \frac{1}{s} F(s) \quad (79)$$

- Differentiation:

$$\frac{d}{dt} f(t) \Leftrightarrow sF(s) \quad (80)$$

- Shift:

$$f(t - T) \Leftrightarrow e^{-Ts} F(s) \quad (81)$$

Proof:

$$\mathbf{LT}(h(t - T)) = \int_0^\infty \underbrace{h(t - T)}_{\text{causal}} e^{-st} dt \quad (82)$$

$$= \int_0^\infty h(t) e^{-s(t+T)} dt \quad (83)$$

$$= \int_0^\infty h(t) e^{-st} e^{-sT} dt \quad (84)$$

$$= e^{-sT} \underbrace{\int_0^\infty h(t) e^{-st} dt}_{H(s)} \quad (85)$$

$$= e^{-sT} H(s) \quad (86)$$

- Convolution:

$$f(t) * g(t) \Leftrightarrow F(s)G(s) \quad (87)$$

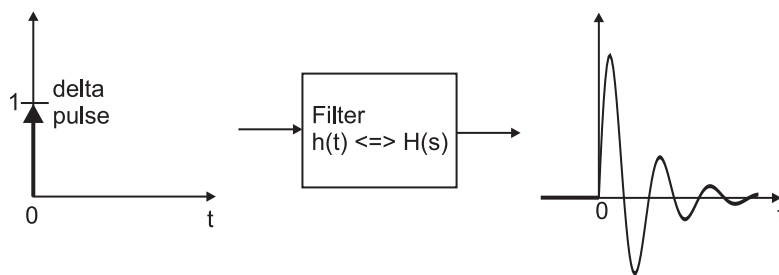


Figure 14: General idea of a filter and how to describe it: either with its impulse response or with it's Laplace transforms.

8.4 Filters

Fig. 14 presents a causal filter as a black box. We send in a signal and we get a signal out of it. The operation of the filter is that of a convolution of the input signal or a multiplication in the Laplace space:

$$g(t) = h(t) * x(t) \Leftrightarrow Y(s) = H(s) \cdot X(s) \quad (88)$$

8.4.1 How to characterise filters?

1. Impulse Response

$$x(t) = \delta(t) \quad \leftarrow \text{delta pulse} \quad (89)$$

$$h(t) = y(t) \quad \leftarrow \text{impulse response} \quad (90)$$

The filter is fully characterised by its impulse response $h(t)$

2. **Transfer function** The Laplace transform of the impulse response is called *Transfer Function*. With the argument $j\omega$ we get the frequency response of the filter. What does the frequency response tell us about the filter? The absolute value of the

$$|H(j\omega)| \quad (91)$$

gives us the amplitude or magnitude for every frequency (compare the Fourier transform). The angle of the term $H(j\omega)$ gives us the phase shift:

$$\phi = \arg(H(j\omega)) \quad (92)$$

of the filter. In this context the *group delay* can be defined as:

$$\tau_\omega = \frac{d\phi(\omega)}{d\omega} \quad (93)$$

which is delay for a certain frequency ω . In many applications this should be kept constant for all frequencies.

8.5 The z-transform

The Laplace transform is for continuous causal signals but in DSP we have sampled signals. So, we need to investigate what happens if we feed a sampled signal:

$$x(t) = \sum_{n=0}^{\infty} x(n)\delta(t - nT) \quad \text{Sampled signal} \quad (94)$$

into the Laplace transform:

$$X(s) = \sum_{n=0}^{\infty} x(n)e^{-snT} \quad (95)$$

$$= \int_0^{\infty} x(n)e^{-st} dt \quad (96)$$

$$= \sum_{n=0}^{\infty} x(n) \underbrace{(e^{-sT})^n}_{z^{-1}=e^{-sT}} \quad (97)$$

$$= \sum_{n=0}^{\infty} x(n)(z^{-1})^n \quad \text{z-transform} \quad (98)$$

What is $e^{-sT} = z^{-1}$? It's a unit delay (see Eq. 81).

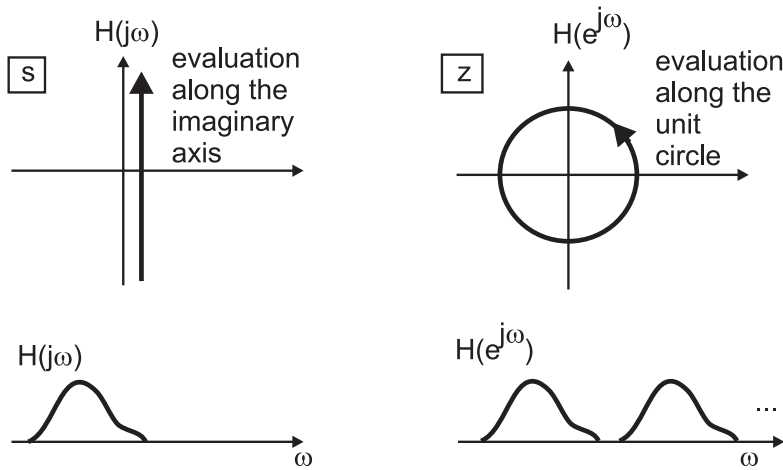


Figure 15: Comparing the calculation of the frequency response in the continuous case and sampled case.

8.6 Frequency response of a sampled filter

Reminder: for analogue Signals we had $H(s)$ with $s = j\omega$. Let's substitute s by z : $z^{-1} = e^{-sT}$ which gives us in general a mapping between the continuous domain and the sampled domain:

$$z = e^{sT} \quad (99)$$

With $s = j\omega$ this gives us $z = e^{j\omega T}$ and therefore the frequency response of the digital filter is:

$$H(e^{j\omega T}) \quad (100)$$

which then leads to the amplitude and phase response:

- Amplitude/Magnitude:

$$|H(e^{i\omega})| \quad (101)$$

- Phase

$$\arg H(e^{i\omega}) \quad (102)$$

Fig. 15 shows the difference between continuous and sampled signal. While in the continuous case the frequency response is evaluated along the imaginary axis, in the sampled case it happens along the unit circle which makes the response periodic! This is a subtle implication of the sampling theorem.

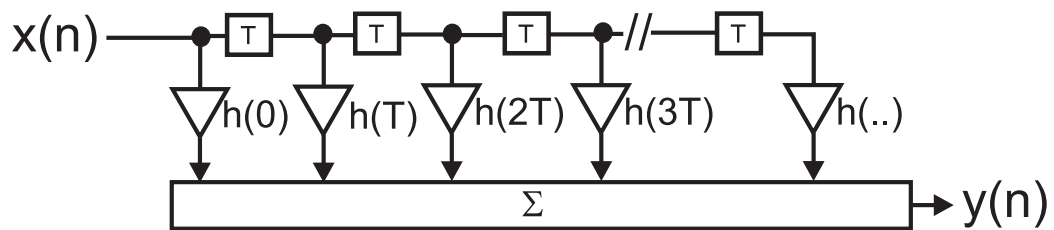


Figure 16: FIR filter using the impulse response of the analogue filter $h(t)$

8.7 FIR Filter

What happens if we sample the impulse response $h(t)$ of an analogue filter? Let's find out:

$$h(t) = \sum_{n=0}^{\infty} h(nT)\delta(t - nT) \quad (103)$$

If we transform it to the Laplace space it looks like this:

$$H(s) = \sum_{n=0}^{\infty} h(nT) \underbrace{(e^{-sT})^n}_{z^{-1}} \quad (104)$$

Remember that e^{-sT} has a very handy meaning: it is a delay by the unit time step (Eq. 81). Thus z^{-n} is a delay by n time steps. We can rewrite Eq. 104:

$$H(z) = \sum_{n=0}^{\infty} h(nT)(z^{-1})^n \quad (105)$$

This is the z-transform of the impulse response $h(t)$ of the filter.

We filter now the signal $X(z)$ with $H(z)$:

$$H(z)X(z) = \underbrace{\sum_{n=0}^{\infty} h(nT)z^{-n}}_{H(z)} X(z) \quad (106)$$

This sum is a direct recipe how to filter the signal $X(z)$. We only need the impulse response of the filter $h(nT)$ and we can set up a digital filter (see Fig. 16). Of course in practise this impulse response cannot run till infinity but only for a limited number of samples. These are often called “taps”. So for example a filter with 100 samples of $h(nT)$ has 100 “taps”. This in turn then requires a delay line which can hold $M = 100$ samples:

$$H(z)X(z) = \underbrace{\sum_{m=0}^M h(mT)z^{-m}}_{H(z)} X(z) \quad (107)$$

This is the formula of an Finite Impulse Response filter where we have sampled an analogue impulse response $h(t)$ at time intervals of T . However, usually the impulse response of the filter is directly derived in the digital domain where the argument n of $h(n)$ represents just the sample index n and the sampling interval is implicitly the inverse of the sampling rate. For that reason one needs to distinguish between:



Figure 17: Digital FIR filter with normalised frequency

$$H(z)X(z) = \sum_{m=0}^M h_{\text{analogue}}(mT)z^{-n} X(z) \quad (108)$$

$$H(z)X(z) = \sum_{m=0}^M h_{\text{digital}}(m)z^{-n} X(z) \quad (109)$$

$$= \sum_{m=0}^M h(m)z^{-m} X(z) \quad (110)$$

Eq 108 is the filter where the impulse response is based on an analogue filter and Eq. 109 is based on the impulse response originating from a digital filter design (see section 8.7.6 and 8.7.7) where the frequencies are normalised frequencies from 0 to 0.5. From now on we will always refer to the “digital” one (Eq. 109) and the subscript will be omitted (Eq. 110 & Fig. 17).

In the time domain the FIR filter is then doing the following operation:

$$y(n) = \sum_{m=0}^{M-1} h(m)x(n-m) \quad (111)$$

8.7.1 Generally FIR filters do not perform a convolution operation

Remember that we have derived the FIR filters from the analogue domain by filtering a signal (Eq.106). The timedomain representation of this is a discrete and causal convolution operation (using the sampled Laplace transform). However, the problem is that this runs to *infinite time* (see Eq 111) and would require an infinite number of delay steps and thus infinite time to complete. So the FIR filter at best performs an approximation of the convolution operation but there are serious problems with the finite number of taps (N) which requires a technique called windowing which moves it even further away from a convolution operation. On the other hand if the impulse response is strictly limited in time

one can use this property for example for matched filters. However, the bottom-line is that sweeping statements such as that “FIR filters perform a convolution” are *generally wrong* if nothing is known of the impulse response.

8.7.2 FIR filter implementations

No matter the design process the implementation is always the same: we need a delay line for the incoming signal and then weight the delayed outputs by the different coefficients and for that reason are called (water-) “taps” (see Eqs. 108 & 109). We are now presenting different implementations.

- C++: This is a simple example of a filter which stores the values in a simple linear buffer `bufferFIR` which stores the delayed values. The coefficients are stored in `coeffFIR`.

```
float filter(float value) {
    // shift
    for (int i=taps-1;i>0;i--) {
        bufferFIR[i]=bufferFIR[i-1];
    }
    //store new value
    bufferFIR[0]=value;
    //calculate result
    for (int i=0;i<taps;i++) {
        output +=bufferFIR[i]*coeffFIR[i];
    }
    return output;
}
```

- Python: Here, the FIR filter is implemented as a class which receives the FIR filter coefficients in the constructor and then filters a signal sample by sample in the function `filter`:

```
class FIR_filter:
    def __init__(self,_coefficients):
        self.ntaps = len(_coefficients)
        self.coefficients = _coefficients
        self.buffer = np.zeros(self.ntaps)

    def filter(self,v):
        self.buffer = np.roll(self.buffer,1)
```

```

self.buffer[0] = v
return np.inner(self.buffer,self.coefficients)

```

which again processes the signal sample by sample. It uses the numpy “roll” command to shift the samples and then the inner product to calculate the weighted sum between the buffer and the coefficients.

If one wants to filter a whole array one can use Python’s lfilter command:

```

import scipy.signal as signal
y = signal.lfilter(h,1,x)

```

This filters the signal x with the impulse response h . Note that this operation is on an array and thus a-causal.

More sophisticated code can be found in [Teukolsky et al. \(2007\)](#). This book is strongly recommended for any C programmer who needs efficient solutions.

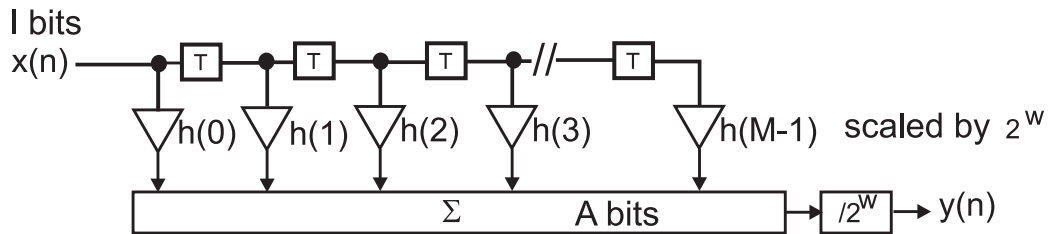


Figure 18: Fixed point FIR filter. The output signal is bit shifted to the right by w bits while the coefficients are scaled up by 2^w .

8.7.3 Fixed point FIR filters

These filters receive integer numbers as input, perform integer multiplications/additions and their outputs are integer as well. Thus, these filters do not require a floating point unit on a processor.

Fig. 18 shows a fixed point FIR filter. The input $x(n)$ is an integer variable with I bit integers, the accumulator is an integer variable with A bits and the output as well (usually the same as the input in terms of bit width).

In contrast to a floating point FIR filter we need to scale up the coefficients so that they use full the integer range to avoid quantisation errors. For example if the coefficients of $h(n)$ range from -0.75 and $+0.75$ and we have signed 16 bit integers then the scaling factor is 2^W , $W = 15$.

However, the accumulator A which collects the data needs to have more bits because it receives scaled input values at I bits precision and these multiplied by factor 2^W . If we have M taps then the additional bits we need is $\log_2 M$. The total number of bits we need in the accumulator in the worst case are:

$$A = I + W + \log_2 M \quad (112)$$

However, this is the worst case scenario because if the gain of the FIR filter is below one then the summations by the M taps will only create *temporary* overflows because integer numbers are cyclic in their representation. In case the gain of the FIR filter is below one this can be relaxed:

$$A = I + W \quad (113)$$

The actual bitwidth of the accumulator is usually the next integer size available and also makes sure that in case the gain goes slightly over one in an unexpected case that the filter still works. For example if I has 16 bits the accumulator has probably 32 bits.

8.7.4 Constant group delay or linear phase filter

So far the FIR filter has no constant group delay which is defined by:

$$\tau_\omega = \frac{d\phi(\omega)}{d\omega} \quad (114)$$

This means that different frequencies arrive at the output of the filter earlier or later. This is not desirable. The group delay τ_ω should be constant for all frequencies ω so that all frequencies arrive at the same time at the output y of the filter.

A constant group delay can be achieved by restricting ourselves to the transfer function:

$$H(e^{i\omega}) = B(\omega)e^{-i\omega\tau+i\phi} \quad (115)$$

where $B(\omega)$ is real and the phase is only defined by the exponential. The phase of this transfer function is then trivially $\omega\tau + \phi$. The group delay is the derivative of this term which is constant.

Eq. 115 now imposes restrictions on the impulse response $h(t)$ of the filter. To show this we use the inverse Fourier transform of Eq. 115 to get the impulse response:

$$h(n) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} H(e^{i\omega})e^{i\omega n} d\omega \quad (116)$$

After some transformations we get:

$$h(n + \tau) = \frac{1}{2\pi} e^{i\phi} b(n) \quad (117)$$

where $b(n)$ represents the Fourier coefficients of $B(\omega)$. Since B is real the coefficients $b(n)$ have the property $b(n) = b^*(-n)$. Thus we get a second impulse response:

$$h(n + \tau) = \frac{1}{2\pi} e^{i\phi} b^*(-n) \quad (118)$$

Now we can eliminate b by equating Eq. 117 and Eq. 118 which yields:

$$h(n + \tau) = e^{2i\phi} n^*(-n + \tau) \quad (119)$$

With the shift τ we have the chance to make the filter “more” causal. We can shift the impulse response in positive time to get $h(n)$ zero for $n < 0$. In a practical application we shift the impulse response by half the number of delays. If we have a filter with M taps we have to delay by $\tau = M/2$.

The factor $e^{2i\phi}$ restricts the values of ϕ because the impulse response must be real. This gives us the final FIR design rule:

$$h(n + M/2) = (-1)^k h(-n + M/2) \quad (120)$$

where $k = 0$ or 1 . This means that the filter is either symmetric or antisymmetric and the impulse response has to be delayed by $\tau = M/2$.

8.7.5 Window functions

So far we still have a infinite number of coefficients for the FIR filter because there's no guarantee that the impulse response becomes zero after $M/2$ delays. Thus, we have to find a way to truncate the response without distorting the filter response.

The standard technique is to multiply the coefficients with a window function which becomes and stays zero at a certain coefficient $n > N$ so that Eq. 107 need not to run to infinity:

$$H(z)X(z) = \sum_{n=0}^N \underbrace{h(nT)w(nT)} z^{-n} X(z) \quad (121)$$

1. Rectangular window: truncating the impulse response. Problem: we get ripples in the frequency-response. The stop-band damping is poor
2. Triangular or Bartlett window: greatly improved stop-band attenuation.
3. Hanning and Hamming window: standard windows in many applications.

$$w(n) = \alpha - (1 - \alpha) \cos\left(\frac{2\pi n}{M}\right) \quad (122)$$

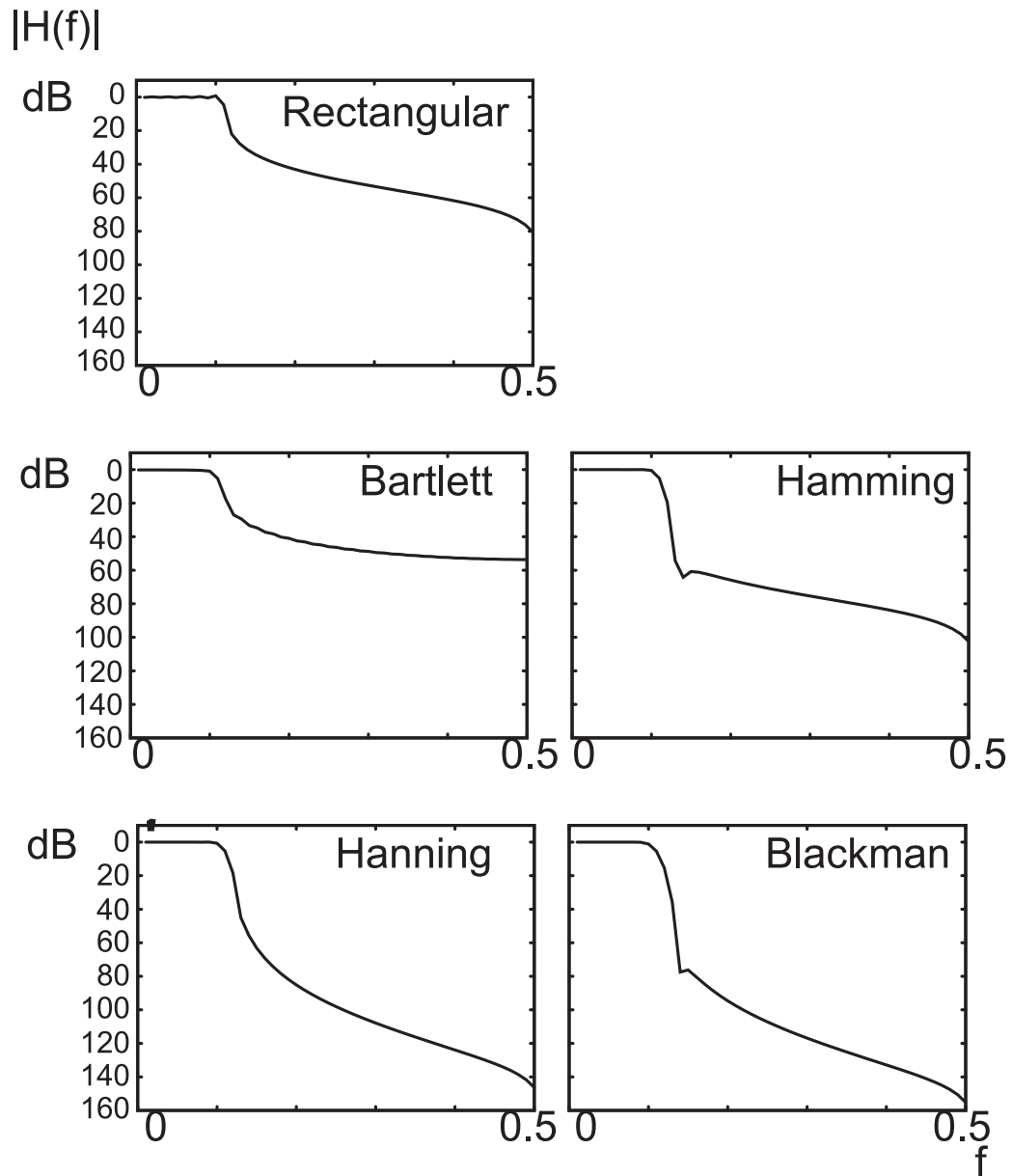


Figure 19: Different window functions applied to a low pass filter (Eq. 133) with cutoff at $f = 0.1$ and 100 taps.

- Hamming: $\alpha = 0.54$
- Hanning: $\alpha = 0.5$

4. Blackman window:

$$w(n) = 0.42 + 0.5 \cos\left(\frac{2\pi n}{M}\right) + 0.08 \cos\left(\frac{4\pi n}{M}\right) \quad (123)$$

5. Kaiser window: control over stop- and passband. No closed form equation available.

To illustrate how window functions influence the frequency response we have taken an impulse response of a lowpass filter ($f_c = 0.1$) and applied different window functions to it (Fig. 19).

Note that the higher the damping the wider the transition from pass- to stopband. This can be seen when comparing the Blackman window with the Hamming window (Fig. 19). For the lowpass filter this seems to be quite similar. However, for a bandstop filter the wider transition width might lead actually to very poor stopband damping. In such a case a Hamming window might be a better choice.

8.7.6 Python code: impulse response from the inverse DFT - The frequency sampling method

Imagine that we want to remove $50Hz$ from a signal with sampling rate of $1kHz$. We define an FIR filter with 100 taps. The midpoint $N/2 = 50$ corresponds to $500Hz$. The $50Hz$ correspond to index 5.

```
f_resp=np.ones(100)
# note we need to add "+1" to the end because the end of
# the range is not included.
f_resp[4:6+1]=0
f_resp[94:96+1]=0
hc=np.fft.ifft(f_resp)
h=np.real(hc)
# this is from index 0 to index 49 on the left
# and on the right hand side from index 50 to index 99
h_shift[0:50]=h[50:100]
h_shift[50:100]=h[0:50]
h_wind=h_shift*hamming(100)
```

To get a nice symmetric impulse response we need to shift the inverse around 50 samples.

8.7.7 FIR filter design from ideal frequency response – The analytical way

For many cases the impulse response can be calculated analytically. The idea is always the same: define a function with the ideal frequency response

$$|H(e^{j\omega})| = \underbrace{B(e^{j\omega})}_{\text{real}} \quad (124)$$

and perform an inverse Fourier transform to get the impulse response $h(n)$.

We demonstrate this with a lowpass filter:

$$|H(e^{j\omega})| = \begin{cases} 1 & \text{for } |\omega| \leq \omega_c \\ 0 & \text{for } \omega_c < |\omega| \leq \pi \end{cases} \quad (125)$$

Use the inverse Fourier transform to get the impulse response:

$$h(n) = \frac{1}{2\pi} \int_{-\pi}^{\pi} H(e^{j\omega}) e^{j\omega n} d\omega \quad (126)$$

$$= \frac{1}{2\pi} \int_{-\omega_c}^{+\omega_c} e^{j\omega n} d\omega \quad (127)$$

$$= \frac{1}{2\pi} \left[\frac{1}{jn} e^{j\omega n} \right]_{-\omega_c}^{+\omega_c} \quad (128)$$

$$= \frac{1}{2\pi jn} (e^{j\omega_c n} - e^{-j\omega_c n}) \quad (129)$$

With these handy equations:

$$\sin z = \frac{1}{2j} (e^{zj} - e^{-zj}) \quad (130)$$

$$\cos z = \frac{1}{2} (e^{zj} + e^{-zj}) \quad (131)$$

we get for the filter:

$$h(n) = \begin{cases} \frac{1}{\pi n} \sin \omega_c n & \text{for } n \neq 0 \\ \frac{\omega_c}{\pi} & \text{for } n = 0 \end{cases} \quad (132)$$

This response is a-causal! However we know that we can shift the response by any number of samples to make it causal! And we have to window the response to get rid of any remaining negative contribution and to improve the frequency response.

Highpass, bandstop and bandpass can be calculated in exactly the same way and lead to the following ideal filter characteristics:

- **Lowpass** with cutoff frequency $\omega_c = 2\pi f_c$:

$$h(n) = \begin{cases} \frac{\omega_c}{\pi} & \text{for } n = 0 \\ \frac{1}{\pi n} \sin(\omega_c n) & \text{for } n \neq 0 \end{cases} \quad (133)$$

- **Highpass** with the cutoff frequency $\omega_c = 2\pi f_c$:

$$h(n) = \begin{cases} 1 - \frac{\omega_c}{\pi} & \text{for } n = 0 \\ -\frac{1}{\pi n} \sin(\omega_c n) & \text{for } n \neq 0 \end{cases} \quad (134)$$

- **Bandpass** with the passband frequencies $\omega_{1,2} = 2\pi f_{1,2}$:

$$h(n) = \begin{cases} \frac{\omega_2 - \omega_1}{\pi} & \text{for } n = 0 \\ \frac{1}{\pi n} (\sin(\omega_2 n) - \sin(\omega_1 n)) & \text{for } n \neq 0 \end{cases} \quad (135)$$

- **Bandstop** with the notch frequencies $\omega_{1,2} = 2\pi f_{1,2}$:

$$h(n) = \begin{cases} 1 - \frac{\omega_2 - \omega_1}{\pi} & \text{for } n = 0 \\ \frac{1}{\pi n} (\sin(\omega_1 n) - \sin(\omega_2 n)) & \text{for } n \neq 0 \end{cases} \quad (136)$$

See [Diniz \(2002, p.195\)](#) for more impulse responses.

Here is an example code for a bandstop filter which fills the array `h` with the analytically calculated impulse response:

```
f1 = 45.0/fs
f2 = 55.0/fs
n = np.arange(-200,200)
h = (1/(n*np.pi))*(np.sin(f1*2*np.pi*n)-np.sin(f2*2*np.pi*n))
h[200] = 1+(f1*2*np.pi-f2*2*np.pi)/np.pi;
h = h * np.hamming(400)
```

After the function has been calculated it is windowed so that the cut off is smooth. It's not very elegant as it causes a division by zero first and then the coefficient at 200 is fixed. Do it better!

8.7.8 Design steps for FIR filters

Here are the design steps for an M tap FIR filter.

1. Get yourself an impulse response for your filter:
 - (a) Create a frequency response "by hand" just by filling an array with the desired frequency response. Then, perform an inverse Fourier transform (see section [8.7.6](#)).

- (b) Define a frequency response analytically. Do an inverse Fourier transform (see section 8.7.7).
 - (c) Use an analogue circuit, get its impulse response and use these as filter coefficients.
 - (d) Dream up directly an impulse response (for example, averagers, differentiators, etc)
2. Mirror the impulse response (if not already symmetrical)
 3. Window the impulse response from an infinite number of samples to M samples.
 4. Move the impulse response to positive time so that it becomes causal (move it $M/2$ steps to the right).

8.7.9 FIR filter design with Python's high level functions

The “firwin” command generates the impulse response of a filter with m taps and also applies a default window function. For example:

```
from scipy.signal import firwin
h = firwin(m, 2*f)
```

generates a lowpass FIR filter with the normalised frequency f . Note the factor two because in scipy the normalisation is *not* the sampling rate but the *nyquist frequency*. To be compatible with the math here the easiest approach is to multiply the frequency by 2. With this command one can also design high pass, bandstop and bandpass filters. Type `help(firwin)` which provides examples for all filter types and which parameters need to be set.

8.7.10 Signal Detection: Matched filter

How can I detect a certain event in a signal? With a correlator. Definition of a correlator?

$$e(t) = \int_0^t \underbrace{s(\tau)}_{\text{signal}} \underbrace{r(\tau)}_{\text{template}} d(\tau) \quad (137)$$

How to build a correlator with a filter? Definition of filtering?

$$e(t) = \int_0^\infty s(\tau) h(t - \tau) d\tau \quad (138)$$

NB. h - integration runs backwards. However we used an int forward! $h(t) := r(T - \tau)$, only valid for $0 \dots T$.

$$e(t) = \int_0^T s(\tau) r(T - (t - \tau)) d\tau \quad (139)$$

$$= \int_0^T s(\tau) r(T - t + \tau) d\tau \quad (140)$$

for $t := T$ we get:

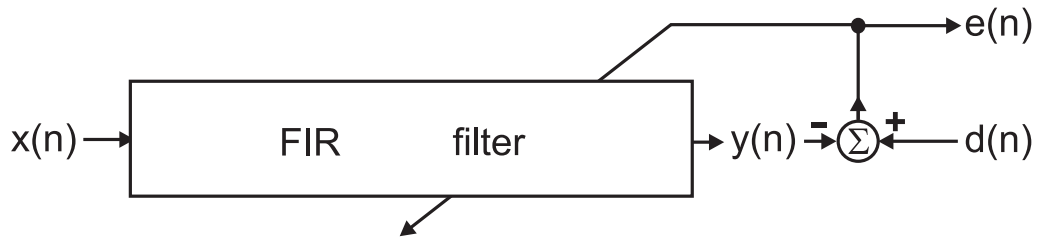
$$e(T) = \int_0^\infty s(\tau) r(\tau) d\tau \quad (141)$$

$$\underbrace{h(t) := r(T - t)}_{\text{matched filter!}} \quad (142)$$

In order to design a detector we just create an impulse response h by reversing the template r in time and constructing an FIR filter with it.

How to improve the matching process? Square the output of the filter!

A



B

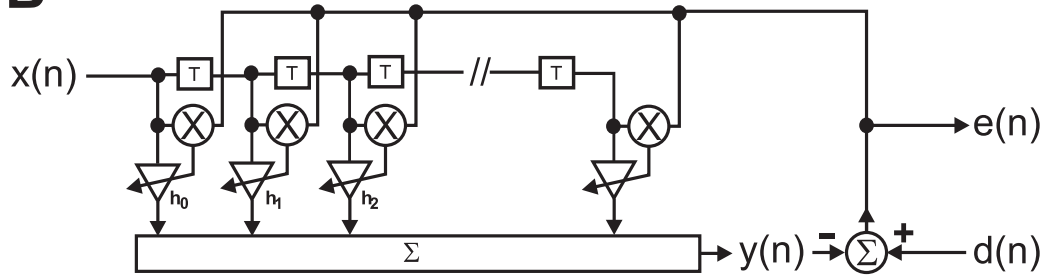


Figure 20: Adaptive FIR filter : A) Negative feedback loop where a desired signal $d(n)$ is compared with the output of the FIR filter $y(n)$ and creates the error signal $e(n)$ which then tunes the FIR filter. B) Internal workings of the adaptive FIR filter where the error signal $e(n)$ is multiplied with the delayed input signal $x(n - m)$ and then changes the FIR filter coefficients h_m .

8.7.11 Adaptive FIR Least Mean Squares (LMS) filters

So far the coefficients of the FIR filter have been constant but now we are going to allow them to change *while* the FIR filter is operating so that it can learn its own coefficients. The idea is to use an error signal $e(n)$ to tune the coefficients $h(n)$ of the FIR filter (see Fig. 20A) by comparing its output $y(n)$ with a desired output $d(n)$:

$$e(n) = d(n) - y(n) \quad (143)$$

and tuning the FIR filter in a negative feedback loop till the average of $e(n)$ is zero.

We need to derive how the error $e(n)$ can change the FIR filter coefficients so that the error is actually minimised. This can be expressed as a so called “gradient descent” which minimises the squared error $\frac{1}{2}e(n)^2$ because both a positive and a negative error are equally bad:

$$\Delta h_m = -\mu \frac{\partial \left(\frac{1}{2}e(n)^2 \right)}{\partial h_m} \quad (144)$$

where $\Delta h_m(n)$ is the change of the coefficient $h_m(n)$ at every time step: $h_m(n+1) = h_m(n) + \Delta h_m(n)$. $\mu \ll 1$ defines how quickly the coefficients h_m change at every time step and is called the “learning rate”. Note here the change in notation of the FIR filter coefficients h_m which are changing much *slower* than the sampled signals $e(n)$, $d(n)$ and $y(n)$ and can still be seen as constant for the realtime filtering operation. Thus, we have n for the sample by sample processing as before and the index m of h_m for the very slowly changing FIR filter coefficients. To gain some intuition why Eq. 144 minimises the error $e(n)$ we look at what happens if we increase the FIR coefficient h_m a little bit, for example, by the small amount of $|\epsilon| \ll 1$: $h_m := h_m + \epsilon$. Then we can observe two cases:

1. The squared error $e(n)^2$ increases so we need to decrease h_m as it makes it worse.
2. The squared error $e(n)^2$ decreases so we need to increase h_m as it makes it better.

This also works in the same way if we decrease h_m by a small amount. Thus from both directions this always minimises the error. This is basically a carrot & stick approach where the coefficients h_m are being “rewarded” if they minimise the error and “punished” if they make the error larger. This approach is called “Least Mean Squares” (LMS) as it minimises the squared error on average. It’s also known from neural networks where the h_m are called the “weights” of a neuron (see delta rule).

Eq 144 now needs to be turned into a form which can directly run in software by solving its partial derivative by inserting Eqs. 143 and 111.

$$\Delta h_m = -\mu \frac{1}{2} \frac{\partial (d(n) - y(n))^2}{\partial h_m} \quad (145)$$

$$= -\mu \frac{1}{2} \frac{\partial \left(d(n) - \sum_{m=0}^{M-1} h_m \cdot x(n-m) \right)^2}{\partial h_m} \quad (146)$$

$$= \mu (d(n) - y(n)) \cdot x(n-m) \quad (147)$$

$$= \mu \cdot e(n) \cdot x(n-m) \quad (148)$$

Note that $x(n-m)$ emerges because of the chain rule when partially differentiating the output $y(n)$ which depends on the sum of the delay lines of the FIR filter (see Eq 111). Eq. 148 is now our “learning” rule which can simply be applied to the FIR filter as shown in Fig 20B.

Now we have the recipe of an adaptive filter where the FIR filter minimises its own error $e(n)$ in a negative feedback loop and generating an output $y(n)$ which is as closely as possible to the desired signal $d(n)$. The signal $y(n)$ is often called “the remover” as it cancels out the signal $d(n)$. However, not everything in $d(n)$ is cancelled out but *only* what is correlated with the input $x(n)$. This means that the error signal $e(n)$ will not be zero but will contain any signal components which are *not* correlated with $x(n)$. This means that the error signal actually is also the cleaned up output signal of the adaptive filter.

Imagine you want to remove 50 Hz mains from a signal

$$d(n) = \text{signal}(n) + 50\text{Hz noise}(n) \quad (149)$$

then one would provide 50 Hz mains via

$$x(n) = 50\text{Hz noise}(n) \quad (150)$$

so that then the filter learns to remove the powerline interference. One might think that the filter will remove everything from $d(n)$ because it will try to minimise the error $e(n)$ but only the 50 Hz in the contaminated signal will be removed because only signal components which are correlated between $x(n)$ and $d(n)$ will be removed. The cleaned up signal will be identical to the error signal:

$$e(n) = \text{signal}(n) \quad (151)$$

8.8 IIR Filter

IIR filter stands for Infinite Impulse Response. Such filters have feedback so that signals are processed in a recursive way. We will see that impulse responses from exponentials can easily be implemented as a recursive filter.

8.8.1 Introduction

As an instructional example we take a very simple analogue filter and sample it. This can then be generalised to more complex situations.

We define a first order filter which can be implemented, for example, as a simple RC network:

$$h(t) = e^{-bt} \quad (152)$$

where its Laplace transform is a simple fraction:

$$H(s) = \frac{1}{s + b} \quad (153)$$

Now we sample Eq. 152:

$$h(t) = \sum_{n=0}^{\infty} e^{-bnT} \cdot \delta(t - nT) \quad (154)$$

and perform a Laplace transform on it:

$$H(s) = \sum_{n=0}^{\infty} e^{-bnT} \underbrace{e^{-nsT}}_{z^{-1n}} \quad (155)$$

which turns into a z-transform:

$$H(z) = \sum_{n=0}^{\infty} e^{-bnT} z^{-n} \quad (156)$$

$$= \sum_{n=0}^{\infty} \left(e^{-bT} z^{-1} \right)^n \quad (157)$$

$$= \frac{1}{1 - e^{-bT} z^{-1}} \quad (158)$$

Consequently the analogue transfer function $H(s)$ transfers into $H(z)$ with the following recipe:

$$H(s) = \frac{1}{s + b} \quad \Leftrightarrow \quad H(z) = \frac{1}{1 - e^{-bT} z^{-1}} \quad (159)$$

Thus, if you have the poles of an analogue system and you want to have the poles in the z-plane you can transfer them with:

$$z_{\infty} = e^{s_{\infty}T} \quad (160)$$

This also gives us a stability criterion. In the Laplace domain the poles have to be in the left half plane (real value negative). This means that in the sampled domain the poles have to lie within the unit circle.

The same rule can be applied to zeros

$$z_0 = e^{s_0T} \quad (161)$$

together with Eq. 160 this is called “The matched z-transform method”. For example $H(s) = s$ turns into $H(z) = 1 - z^{-1}e^{0T} = 1 - z^{-1}$ which is basically a DC filter.

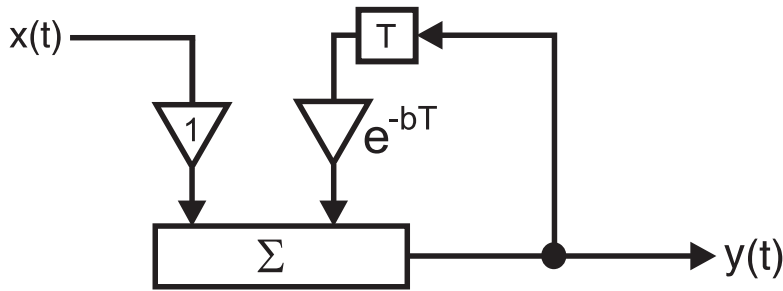


Figure 21: IIR filter

8.8.2 Determining the data-flow diagram of an IIR filter

To get a practical implementation of Eq. 158 we have to see it with its input- and output-signals:

$$Y(z) = X(z)H(z) \quad (162)$$

$$= X(z) \frac{1}{1 - e^{-bT}z^{-1}} \quad (163)$$

$$= Y(z)z^{-1}e^{-bT} + X(z) \quad (164)$$

We have to recall that z^{-1} is the delay by T . With that information we directly have a difference equation in the temporal domain:

$$y(nT) = y([n - 1]T)e^{-bT} + x(nT) \quad (165)$$

This means that the output signal $y(nT)$ is calculated by adding the weighted and delayed output signal $y([n - 1]T)$ to the input signal $x(nT)$. How this actually works is shown in Fig. 21. The attention is drawn to the *sign inversion* of the weighting factor e^{-bT} in contrast to the transfer function Eq. 159 where it is $-e^{-bT}$. In general the recursive coefficients change sign when they are taken from the transfer function.

8.8.3 General form of IIR filters

A filter with forward and feedback components can be represented as:

$$H(z) = \frac{\sum_{k=0}^r B_k z^{-k}}{1 + \sum_{l=1}^m A_l z^{-l}} \quad (166)$$

where B_k are the FIR coefficients and A_l the recursive coefficients. Note the signs of the recursive coefficients are *inverted* in the actual implementation of the filter. This can be seen when the function $H(z)$ is actually multiplied with an input signal to obtain the output signal (see Eq. 165 and Fig. 21). The “1” in the denominator represents actually the output of the filter. If this factor is not one then the output will be scaled by that factor. However, usually this is kept one.

In Python filtering is performed with the command:

```
import scipy.signal as signal
Y = signal.lfilter(B,A,X)
```

where B are the FIR coefficients, A the IIR coefficients and X is the input. For a pure FIR filter we just have:

```
Y = signal.lfilter(B,1,X)
```

The “1” represents the output.

8.8.4 IIR filter topologies

In real applications one would create a *chain of 2nd order IIR filters*. This has numerous advantages:

1. The optimisation can be done within these simple structures, for example omitting array operations completely and coding the 2nd order filters in assembly while keeping the higher level operations in C or C++.



Figure 22: A) Direct Form I filter which one accumulator and B) Direct Form II filter with two accumulators.

2. Control of stability issues: high order recursive systems might become unstable and it is very hard to predict when this might happen. On the other hand a chain of 2nd order filters can be made stable with ease. One can focus on stability in 2nd order systems which is well understood and manageable.
3. Complex conjugate pole pairs are generated naturally in analogue filter design and directly translate to 2nd order IIR structures. Analogue design is usually described as 2nd order structures (=LCR) so that any transform from analogue to digital just needs to be of 2nd order!

Fig. 22 shows the most popular filter topologies: Direct Form I and II. Because of the linear operation of the filter one is allowed to do the FIR and IIR operations in different orders. In Direct Form I we have one accumulator and two delay lines whereas in the Direct Form II we have two accumulators and one delay line. Only the Direct Form I is suitable for integer operations.

A Python class of a direct form II filter can be implemented with a few lines:

```
class IIR_filter:
    def __init__(self, _num, _den):
        self.numerator = _num
        self.denominator = _den
        self.buffer1 = 0
        self.buffer2 = 0

    def filter(self, v):
        input=0.0
        output=0.0
        input=v
```

```

output=(self.numerator[1]*self.buffer1)
input=input-(self.denominator[1]*self.buffer1)
output=output+(self.numerator[2]*self.buffer2)
input=input-(self.denominator[2]*self.buffer2)
output=output+input*self.numerator[0]
self.buffer2=self.buffer1
self.buffer1=input
return output

```

Here, the two delay steps are represented by two variables `buffer1` and `buffer2`.

In order to achieve higher order filters one can then just chain these 2nd order filters. In Python this can be achieved by storing these in an array of instances of this class.

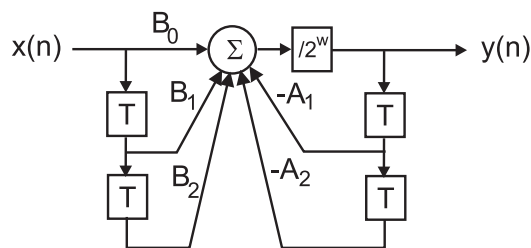


Figure 23: Direct Form I filter with fixed point arithmetic.

8.8.5 Fixed point IIR filters

Fig. 23 shows the implementation of a fixed point IIR filter. It's a Direct Form I filter with one accumulator so that temporary overflows can be compensated. The original floating point IIR coefficients are scaled up by factor 2^w so that they max out the range of the integer coefficients. After the addition operation they are shifted back by w bits to their original values.

For example if the largest IIR coefficient is 1.9 and we use 16 bit signed numbers (max number is 32767) then one could multiply all coefficients with 2^{14} .

Then one needs to assess the maximum value in the accumulator which is much more difficult than for FIR filters. For example, resonators can generate high output values with even small input values so that it's advisable to have a large overhead in the accumulator. For example if the input signal is 16 bit and the scaling factor is 14 bits then the signal will certainly occupy $16+14 = 30$ bits. With an accumulator of 32 bits that gives only a headroom of 2 bits so the output can only be 4 times larger than the input. A 64 bit accumulator is a safe bet.

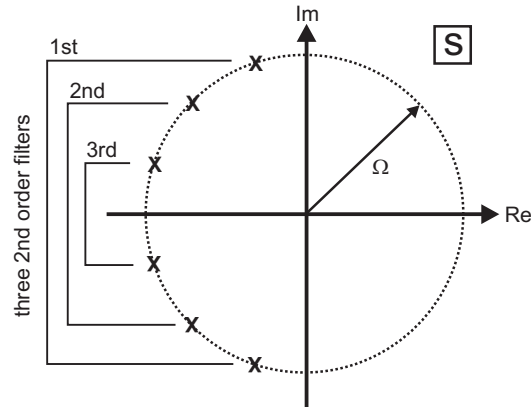


Figure 24: Pole placements of the Butterworth filter. Its analogue cutoff frequency is $\Omega = 2\pi F$ where F is the analogue cutoff frequency in Hertz. This filter can be implemented as a chain of 2nd order IIR filters (biquads) by using the complex conjugate pairs for the different 2nd order IIR filters.

8.8.6 Filter design based on analogue filters

Most IIR filters are designed from analogue filters by transforming from the continuous domain $(h(t), H(s))$ to the sampled domain $(h(n), H(z))$.

A list of popular analogue transfer functions being used for IIR filter design:

- **Butterworth:** All poles lie on the left half plane equally distributed on a half circle with radius $\Omega = 2\pi F$ which is shown in Fig. 24. The Butterworth filter is by far the most popular filter. Here are its properties:
 - monotonic frequency response
 - only poles, no zeros
 - the Poles have analytical solution and very easy to calculate
 - no constant group delay but usually acceptable deviation from a strict constant group delay for many applications.
- **Chebyshev Filters:**

$$|H(\Omega)|^2 = \frac{1}{1 - \varepsilon^2 T_N^2(\Omega/\Omega_p)} \quad (167)$$

where T = Chebyshev polynomials. These filters have either ripples in the stop or passband depending on the choice of polynomials. As with

Butterworth the polynomials have analytical solutions for the poles and zeros of the filter so that their design again is straightforward.

- **Bessel Filter:**

- Constant Group Delay
- Shallow transition from stop- to passband
- No analytical solution for the poles/zeros.

How to transform the analogue lowpass filters into highpass, bandstop or bandpass filters? All the analogue transfer functions above are lowpass filters. However this is not a limitation because there are handy transforms which take an analogue lowpass transfer function or poles/zeros and then transforms them into highpass, bandpass and stopband filters. These rules can be found in any good analogue filter design book and industry handouts, for example from Analog Devices.

Bilinear Transform: transforming the analogue transfer function into a digital one: As a next step these analogue transfer functions $H(s)$ need to be transformed to digital transfer functions $H(z)$. This could be done by the matched z-transform (Eq. 161). However, the problem with these methods is that they map frequencies 1:1 between the digital and analogue domain. Remember: in the sampled domain there is no infinite frequency but $F_s/2$ which is the Nyquist frequency. This means that we never get more damping than at Nyquist $F_s/2$. This is especially a problem for lowpass filters where damping increases the higher the frequency.

The solution is to map **all** analogue frequencies from $0 \dots \infty$ to the sampled frequencies $0 \dots 0.5$ in a non-linear way:

$$-\infty < \Omega < \infty \Rightarrow -\pi \leq \omega \leq \pi \quad (168)$$

This is called *Bilinear Transformation*:

$$s = \frac{2}{T} \frac{z - 1}{z + 1} \quad (169)$$

This rule replaces all s with z in our analogue transfer function so that it's now digital. However, the cutoff frequency Ω_c is still an analogue one but if we want to design a digital filter we want to specify a digital cutoff frequency. Plus remember that the frequency mapping is non-linear so that there is non-linear mapping also of the cutoff.

In the analogue domain the frequency is given as $s = j\Omega$ and in the sampled domain as $z = e^{j\omega}$. With the definition of the bilinear transform we can establish how to map from our desired digital cutoff to the analogue one:

$$j\Omega = \frac{2}{T} \left[\frac{e^{j\omega} - 1}{e^{j\omega} + 1} \right] = \frac{2}{T} j \tan \frac{\omega}{2} \quad (170)$$

This derivation is useful for two purposes. It shows the non-linear mapping between the digital and analogue world and it provides a recipe how to calculate our analogue cutoff frequency.

That the bilinear transform is a **nonlinear** mapping between the analogue world and the digital world can be directly seen by just omitting the j :

$$\Omega = \frac{2}{T} \tan \frac{\omega}{2} \quad (171)$$

This also means that the cut-off frequency of our analogue filter is changed by the bilinear transformation. Consequently, we need to apply the same transformation to the cutoff frequency itself:

$$\Omega_c = \frac{2}{T} \tan \frac{\omega_c}{2} \quad (172)$$

where T is the sampling interval. This is often called “pre-warp” but is simply the application of the same rule to the cut-off frequency as what the bilinear transform does to the transfer function $H(s)$. It has also another important result: we can now finally specify our cut-off in the sampled domain in normalised frequencies $\omega_c = 2\pi f_c$ by using Eq. 172. After all we just use the analogue filter as a vehicle to design a digital filter.

We can now list our design steps.

IIR filter design steps:

1. Choose the cut-off frequency of your digital filter ω_c .
2. Calculate the analogue cutoff frequency $\omega_c \rightarrow \Omega_c$ with Eq. 172
3. Choose your favourite analogue lowpass filter $H(s)$, for example Butterworth.
4. Replace all s in the analogue transfer function $H(s)$ by $s = \frac{2}{T} \frac{z-1}{z+1}$ to obtain the digital filter $H(z)$
5. Change the transfer function $H(z)$ so that it only contains negative powers of z (z^{-1}, z^{-2}, \dots) which can be interpreted as delay lines.

6. Build your IIR filter!

For filter-orders higher than two one needs to develop a different strategy because the bilinear transform is a real *pain* to calculate for anything above the order of two. Nobody wants to transform high order analogue transfer functions $H(s)$ to the $H(z)$ domain. However, there is an important property of all analogue transfer functions: they generate complex conjugate pole pairs (plus one real pole if of odd order) which suggest a chain of 2nd order IIR filters straight away (see Fig. 24). Remember that a complex conjugate pole pair creates a 2nd order IIR filter with two delay steps. A real pole is a 1st order IIR filter with one delay but is often also implemented as a 2nd order filter where the coefficients of the 2nd delay are kept zero.

The design strategy is thus to split up the analogue transfer function $H(s)$ in a chain of 2nd order filters $H(s) = H_1(s)H_2(s)H_3(s)\dots$ and then to apply the bilinear transform on every 2nd order term separately. Using this strategy you only need to calculate the bilinear transform once for a 2nd order system (or if the order is odd then also for a 1st order one) but then there is no need to do any more painful bilinear transforms. This is standard practise in IIR filter design.

Time or frequency domain? The transforms from analogue to digital alter both the temporal response and the frequency response. However the different transforms have different impact. The bilinear transform guarantees that the digital filter uses the whole frequency range of the analogue filter from zero to infinity which is the best solution for frequency domain problems. However, if one needs to reproduce the temporal response an analogue filter as faithfully as possible then the matched z transform is best because it's based (as the name suggests) on the impulse invariance method which aims to preserve the temporal behaviour of the filter (exact for pole-only transfer functions).

Identical frequency response required: Bilinear transform

Identical temporal behaviour required: Matched z-transform

Of course a 100% match won't be achieved but in practise this can be assumed. In most cases the bilinear transform is the transform of choice. The matched z transform can be very useful, for example in robotics where timing is important.

8.8.7 Adaptive IIR filter: The Kalman filter

Often signals are contaminated by high frequency noise where the spectrum of the noise is changing or not known. The cutoff of a fixed low pass filter is not

known or might be changing so that we need to adapt the cut-off continuously. One example is a Kalman filter which maximises the *predictability* of the filtered signal. It's an adaptive lowpass filter which increases the predictability of a signal because it smoothes it.

$$H(z) = \frac{b}{1 - az^{-1}} \quad (173)$$

The parameter a determines the cutoff frequency. Frequency Response :

$$|H(e^{j\omega})| = \left| \frac{b}{1 - ae^{-j\omega}} \right| \quad (174)$$

Let's re-interpret our low-pass filter in the time domain:

$$\underbrace{y(n)}_{\text{actual estimate}} = a(n) \underbrace{y(n-1)}_{\text{previous estimate}} + b(n) \underbrace{x(n)}_{\text{current data sample}} \quad (175)$$

We would like to have the best estimate for $y(n)$

$$p(n) = E[(y(k) - y_{\text{real}}(k))^2] \quad (176)$$

We need to minimise p which gives us equations for a and b which implements a Kalman filter.

8.9 The role of poles and zeros

Transfer functions contain poles and zeros. To gain a deeper understanding of the transfer functions we need to understand how poles and zeros shape the frequency response of $H(z)$. The position of the poles also determines the stability of the filter which is important for real world applications.

We are going to explore the roles of poles and zeros first with an instructional example which leads to a 2nd order bandstop filter.

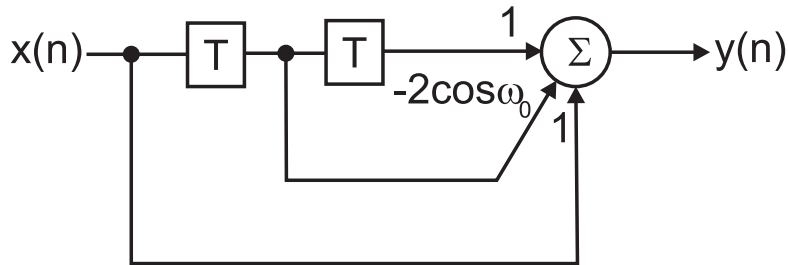


Figure 25: A two tap FIR stopband filter for the frequency ω_0 .

8.9.1 Zeros

Let's look at the transfer function:

$$H(z) = (1 - e^{j\omega_0} z^{-1})(1 - e^{-j\omega_0} z^{-1}) \quad (177)$$

$$= 1 - z^{-1} e^{j\omega_0} + z^{-2} \quad (178)$$

$$= 1 - z^{-1} (e^{j\omega_0} + e^{-j\omega_0}) + z^{-2} \quad (179)$$

$$= 1 - z^{-1} 2 \cos \omega_0 + z^{-2} \quad (180)$$

which leads to the data flow diagram shown in Fig 25.

$$H(e^{j\omega}) = \underbrace{(1 - e^{j\omega_0} e^{-j\omega})(1 - e^{-j\omega_0} e^{-j\omega})}_{2 \text{ zeros}} \quad (181)$$

The zeroes at $e^{j\omega_0}$ and $e^{-j\omega_0}$ eliminate the frequencies ω_0 and $-\omega_0$.

A special case is $\omega_0 = 0$ which gives us:

$$H(z) = (1 - e^0 z^{-1})(1 - e^0 z^{-1}) \quad (182)$$

$$= 1 - 2z^{-1} + z^{-2} \quad (183)$$

a DC filter.

In summary: zeros eliminate frequencies (and change phase). That's the idea of an FIR filter where loads of zeros (or loads of taps) knock out the frequencies in the stopband.

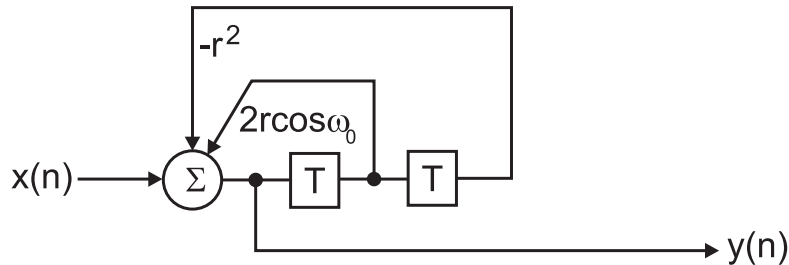


Figure 26: A two tap IIR resonator for the frequency ω_0 and with the amplification r .

8.9.2 Poles

While zeros knock out frequencies, poles amplify frequencies. Let's investigate complex conjugate poles:

$$H(z) = \frac{1}{\underbrace{(1 - r e^{j\omega_0} z^{-1})(1 - r e^{-j\omega_0} z^{-1})}_{2 \text{ poles!}}} \quad (184)$$

which are characterised by their resonance frequency ω_0 and the amplification $0 < r < 1$.

In order to obtain a data flow diagram we need to get powers of z^{-1} because they represent delays.

$$H(z) = \frac{1}{1 - 2r \cos(\omega_0)z^{-1} + r^2z^{-2}} \quad (185)$$

which we multiply with the input signal $Y(z)$:

$$Y(z) = X(z) \frac{1}{1 - 2r \cos(\omega_0)z^{-1} + r^2z^{-2}} \quad (186)$$

$$X(z) = Y(z) - Y(z)2r \cos(\omega_0)z^{-1} + Y(z)r^2z^{-2} \quad (187)$$

$$Y(z) = X(z) + z^{-1}Y(z)2r \cos(\omega_0) - z^{-2}Y(z)r^2 \quad (188)$$

This gives us a second order recursive filter (IIR) which is shown in Fig 26. These complex conjugate poles generate a resonance at $\pm\omega_0$ where the amplitude is determined by r .

8.9.3 Stability

A transfer function $H(z)$ is only stable if the poles lie inside the unit circle. This is equivalent to the analog case where the poles of $H(s)$ have to lie on the left hand side of the complex plane (see Eq. 160). Looking at Eq. 184 it becomes clear that r determines the radius of the two complex conjugate poles. If $r > 1$ then the filter becomes unstable. The same applies to poles on the real axis. Their the real values have to stay within the range $-1 \dots +1$.

In summary: poles generate resonances and amplify frequencies. The amplification is strongest the closer the poles move towards the unit circle. The poles need to stay within the unit circle to guarantee stability.

Note that in real implementations the coefficients of the filters are limited in precision. This means that a filter might work perfectly in python but will fail on a DSP with its limited precision. You can simulate this by forcing a certain datatype in python or you write it properly in C.

8.9.4 Design of an IIR notch filter

Now we combine the filters from the last two sections:

$$H(z) = \frac{1 - 2 \cos(\omega_0)z^{-1} + z^{-2}}{1 - 2r \cos(\omega_0)z^{-1} + r^2z^{-2}} \quad (189)$$

This gives us a notch filter where its width is tunable with the help of $r < 1$. The closer r goes towards 1 the more narrow is the frequency response. This

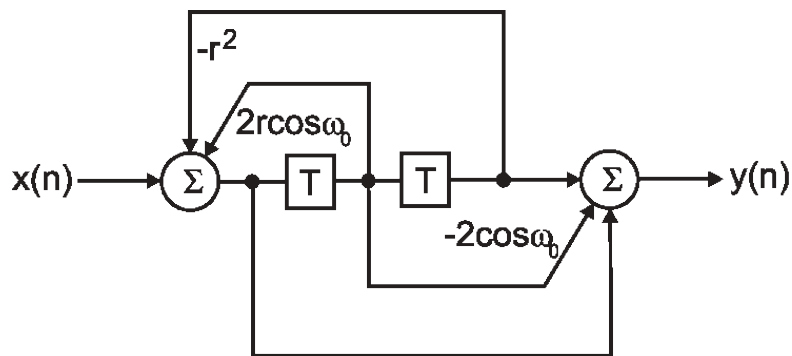


Figure 27: A two tap IIR bandstop filter with tuneable stopband width r for the frequency ω_0 .

filter has two poles and two zeros. The zeros sit on the unit circle and eliminate the frequencies $\pm\omega_0$ while the poles sit within the unit circle and generate a resonance around $\pm\omega_0$. As long as $r < 1$ this resonance will not go towards infinity at $\pm\omega_0$ so that the zeros will always eliminate the frequencies $\pm\omega_0$.

A C++ implementation of such a filter is quite compact and requires no loops:

```
float Iirnotch::filter(float value) {
    float input=0.0;
    float output=0.0;
    // a little bit cryptic but optimized for speed
    input=value;
    output=(numerator[1]*buffer[1]);
    input=input-(denominator[1]*buffer[1]);
    output=output+(numerator[2]*buffer[2]);
    input=input-(denominator[2]*buffer[2]);
    output=output+input*numerator[0];
    buffer[2]=buffer[1];
    buffer[1]=input;
    return output;
}
```

This filter is part of the program comedirecord and is used to filter out 50Hz noise.

8.9.5 Identifying filters from their poles and zeroes

[Proakis and Manolakis \(1996, pp.333\)](#) has an excellent section about this topic

and we refer the reader to have a look. As a rule of thumb a digital lowpass filter has poles where their real parts are positive and a highpass filter has poles with negative real part of the complex poles. In both cases they reside within the unit circle to guarantee stability.

References

- Diniz, P. S. R. (2002). *Digital Signal Processing*. Cambridge university press, Cambridge.
- Proakis, J. G. and Manolakis, D. G. (1996). *Digital Signal Processing*. Prentice-Hall, New Jersey.
- Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (2007). *Numerical Recipes: The Art of Scientific Computing*. Cambridge University Press, 3rd edition.