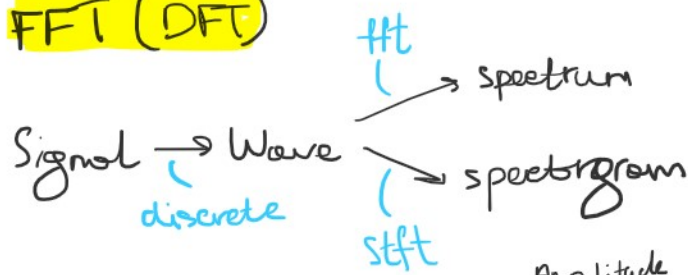
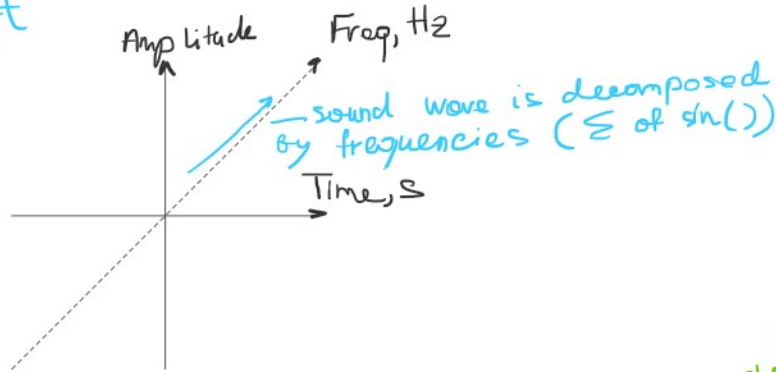


FFT (DFT)



Two domains

- 1) Time domain
- 2) Frequency domain

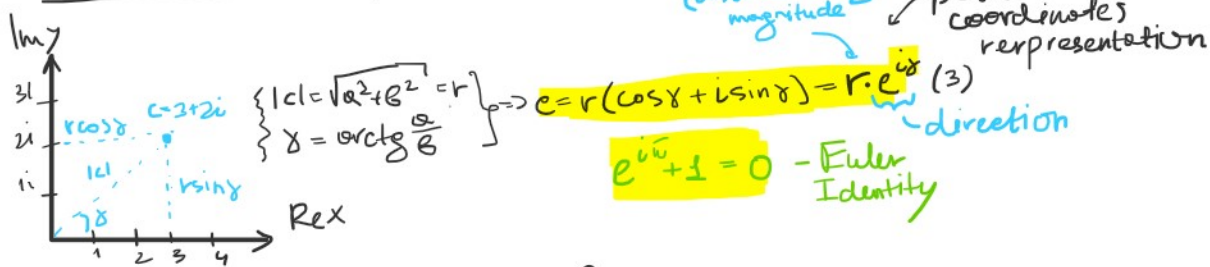


$\sin(2\pi(f \cdot t - \varphi))$ - formula of sin-wave

1 - frequency ($2\pi \cdot f = \omega$)

2 - phase

FT (with complex numbers and math)



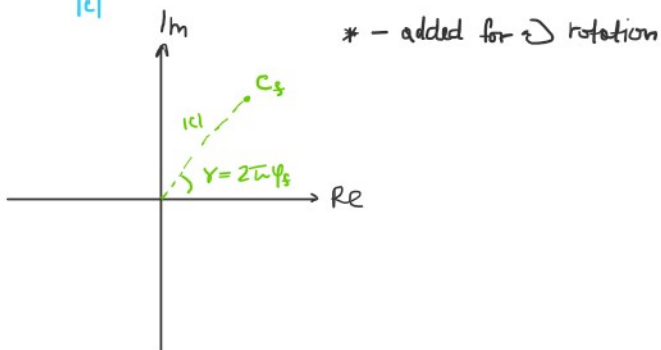
complex numbers recap.

how for FFT.

how to go from (1), (2) to (3)?

Intuition: use magnitude and phase as polar coordinates.

$c_f = \frac{d_f}{|c|} \cdot e^{-i2\pi f t_f}$ - complex FFT coeffs.



$g(t): \mathbb{R} \rightarrow \mathbb{R}$ - original audio signal (time domain)

$\hat{g}(f): \mathbb{R} \rightarrow \mathbb{C}$ - FT, outputs a complex number (c_f)

$\hat{g}(f) = \int_{-\infty}^{\infty} g(t) \cdot e^{-i2\pi f t} dt$

Interpretation of \int - center of mass. The greater $|c_f|$ - the greater the alignment with true sinusoid.

* Since sound is discrete, instead of \int we do \sum over all timestamps.

$d_f = |c_f| \cdot \hat{g}(f) \rightarrow \varphi_f = -\frac{\gamma_f}{2\pi}$

Short-Time Fourier Transform (signal \rightarrow spectrogram)

reminder of DFT: $\hat{x}(k/N) = \sum_{n=0}^{N-1} x(n) \cdot e^{-i2\pi n \cdot \frac{k}{N}}$

Problem: we know what freq. are present, but don't know when. (cause they are averaged over time)

Solution: use windowing

1) $x_w(t) = x(t) \cdot w(t)$ - window function

2) apply fft on window

window size - number of samples $w()$ is applied to

frame size - number of samples in chunk

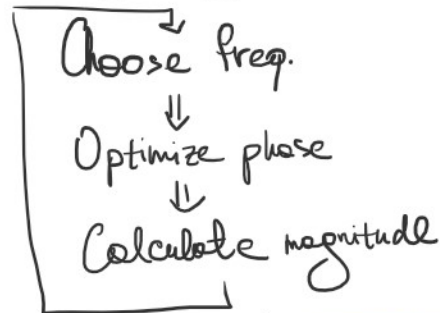
w_s can be $<$ than f_s (then samples outside of w_s will be 0-padded)

1) $x_w(t) = x(t) \cdot w(t)$

2) apply fft on window

Idea of FFT:

- compare signal with different sinusoids
- $\text{FFT}(\text{freq}) \rightarrow \text{magnitude} + \text{phase}$
- high magnitude \Rightarrow sin as signal



$\varphi_f = \arg \max_{\varphi \in [0, 1)} \int s(t) \cdot \sin(2\pi(f t - \varphi)) dt$ (1)

phase $\varphi \in [0, 1)$ signal chosen sinusoid

$d_f = \max_{\varphi \in [0, 1)} \int s(t) \cdot \sin(2\pi(f t - \varphi)) dt$ (2)

magnitude $\varphi \in [0, 1)$

Fourier representation (ift)

$g(t) = \int c_f \cdot e^{i2\pi f t} df$

we add up all weighted sinusoids and get original sound.

DFT case (what is actually going on)

$g(t) \mapsto x(n), n \in \mathbb{N}$

$t = n \cdot T$ (T - sampling rate)

$\hat{x}(f) = \sum_{n=0}^{N-1} x(n) \cdot e^{-i2\pi f n T}$ - discrete version of $\int_{-\infty}^{\infty}$

- consider n where $f_{\text{req}} \neq 0$
- compute for finite set of $\{f_1, \dots, f_M\}$, $M = N$ - number of samples (allows fft and efficient)

$f = \frac{k}{N}$, $k \in [0, M-1] = [0, N-1]$

$F(k) = \frac{k}{NT} = \frac{k \cdot sr}{N}$ - actual freq. value

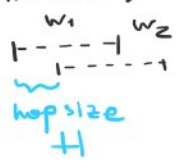
BUT: redundancy - plot will be symmetrical around Nyquist fr. $= \frac{N}{2} > F(N/2) = sr/2$. So we can throw everything after it.

Note on FT to FFT (very fast computing)

- $O(N^2) \rightarrow O(N \log N)$
- $N = 2^x$ ($x \in \mathbb{N}$)

... windows are overlapping. Each window position is defined by hop size parameter.

However, windows are overlapping. Each window position is defined by hop size parameter.



$$S(m, k) = \sum_{n=0}^{N-1} x(n+mH) \cdot w(n) \cdot e^{-i2\pi n \frac{k}{N}}$$

Annotations for the equation:

- m : proxy time value, frame number.
- N : frame size, how many x over frame.
- $x(n+mH)$: all samples from m frame.
- $w(n)$: window function.
- $e^{-i2\pi n \frac{k}{N}}$: phase thingy.

DFT \rightarrow 1-d vector, N complex FT

STFT \rightarrow 2-d matrix, freq-bins \times frames

$$\# \text{freq-bins} = \frac{\text{framesize} + 1}{2} \quad (\text{related to Nyq. freq.})$$

$$\# \text{frames} = \frac{\text{samples} - \text{framesize} + 1}{\text{hopsize}}$$

Parameters study

• $fs = (512, 1024, \dots, 2^n)$

- $fs \downarrow \Rightarrow \text{freq.res.} \uparrow \text{ time.res.} \downarrow$ (time/freq trade off)

$fs \uparrow \Rightarrow \text{freq.res.} \downarrow \text{ time.res.} \uparrow$

• $hs = (256, 512, \dots, 2^m)$, $m = (n-1, n-2, \dots) \sim hs = \frac{fs}{2}$

• $w(n)$ - - better, like Hann window

How to get to spectrogram?

$$Y(m, k) = |S(m, k)|^2$$

\rightarrow real-valued matrix

Example:

signal - 10k samples

$$fs = ws = 1000$$

$$hs = 500$$

$$\text{number of fr. bins} = \frac{1000}{2} + 1 = 501 \rightarrow (0, \frac{fs}{2})$$

$$\text{num. of frames} = (10k - 1k) / 500 + 1 = 19$$

$$\text{STFT} \rightarrow (501, 19)$$