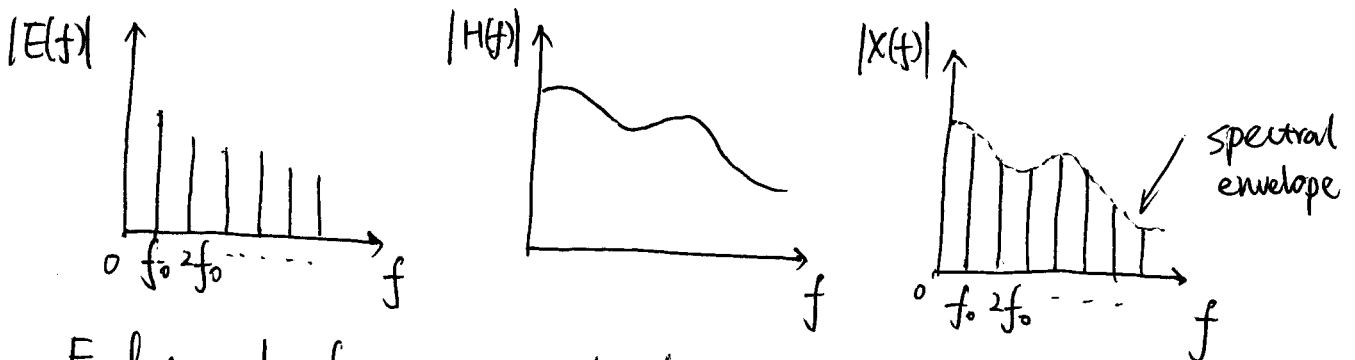Harmonic Sound : e.g. violin, flute, french horn, voice, etc.

Violin :  String vibration $\Rightarrow$ Body resonance $\Rightarrow$ Sound production

(Source / excitation)  (filter)  (Signal)

time domain  $e(t)$  $*$  $h(t)$  $=$  $x(t)$

impulse response
of linear filter

freq. domain  $E(f)$  $\times$  $H(f)$  $=$  $X(f)$



$f_0$ : Fundamental frequency, related to perception "pitch"
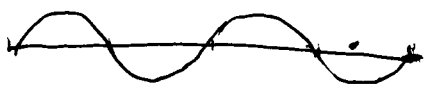
$n f_0$ : Harmonics

String vibration
Standing waves.

a periodic
impulse train



$f_0$

$2 f_0$

$3 f_0$

$4 f_0$

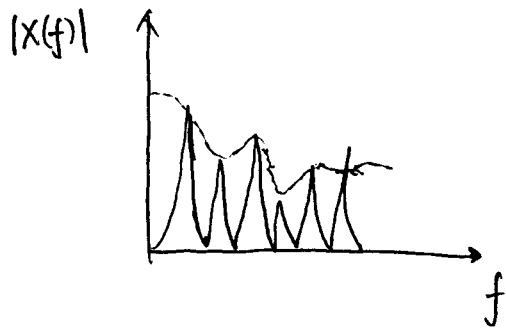Timbre: the thing that people use to discriminate two sounds with the same pitch, loudness, and duration.

What is timbre exactly?

① spectral envelope.

② temporal dynamics.

Question: How do we extract and represent spectral envelope from the final sound production?

First thought: direct calculation from magnitude spectrum

$|X(f)|$

Connect spectral peaks (ideally should represent harmonics) with a smooth curve.

Problem: ① the curve can be sensitive to the peaks. (How do we make sure the peaks are all correct, i.e., represent harmonics?) (What if the signal is not harmonic?)

② Represent the curve non-parametrically? A vector with very high dimension.

Idea: Calculate and represent spectral envelope with a parametric representation.

Linear Predictive Coding.    (Time domain representation of spectral envelope)

Assume current signal is a linear combination of past signals plus the excitation signal.    (Auto-regressive model)

$$x[n] = \sum_{k=1}^{P} a_k \, x[n-k] + e[n]$$

↑ past signal    ↑ excitation (periodic impulse train).

By Z-transform

$$X(z) = \sum_{k=1}^{P} a_k \, X(z) \, z^{-k} + E(z)$$

∴ transfer function $H(z) = \dfrac{X(z)}{E(z)} = \dfrac{1}{1 - \sum_{k=1}^{P} a_k \, z^{-k}}$

(Resonance filter frequency response)

(all-pole model)

We can use $\{a_k\}_{k=1,\cdots,P}$ to represent the filter, i.e., the spectral envelope. $P$ is the model order.

Compared to the signal, the filter changes relatively slowly.

For musical instruments, if it doesn't change basically, as the instrument body doesn't change.

For speech, it changes little within 20 ms.

We perform short-time analysis of the signal to estimate $a_k$.

$$e[n] = x[n] - \sum_{k=1}^{P} a_k x[n-k]$$

Remember $e[n]$ is peoridic impulse train, its value is small most of the time

It's resonable to minimize the mean square of $e[n]$ within a short time

$$\mathcal{E} \triangleq \sum_n e^2[n] = \sum_n \left( x[n] - \sum_{k=1}^{P} a_k x[n-k] \right)^2$$

$$\frac{\partial \mathcal{E}}{\partial a_k} = 2 \sum_n \left( x[n] - \sum_{i=1}^{P} a_i x[n-i] \right) \cdot (-x[n-k]) = 0$$
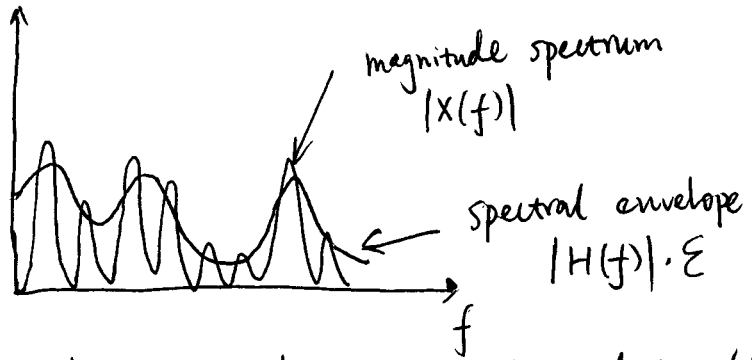
$$\therefore \sum_n x[n] x[n-k] = \sum_{i=1}^{P} a_i \sum_n x[n-i] x[n-k] \quad \text{for } i=1,\cdots,P.$$

P equations, P unknowns.

Let $\varphi[i,k] = \sum_n x[n-i] x[n-k]$    Autocorrelation. $R(|i-k|)$

then we have $\varphi[0,k] = \sum_{i=1}^{P} a_k \varphi[i,k]$  for $k=1,\cdots,P$

frequency response of $H(z)$



magnitude spectrum $|X(f)|$

spectral envelope $|H(f)| \cdot \mathcal{E}$

$\mathcal{E} \cdot |H(f)|$ becomes less smoother and closer to $|X(f)|$ if increase P.

Here we use "energy matching criterion" to define/draw the spectral envelope.
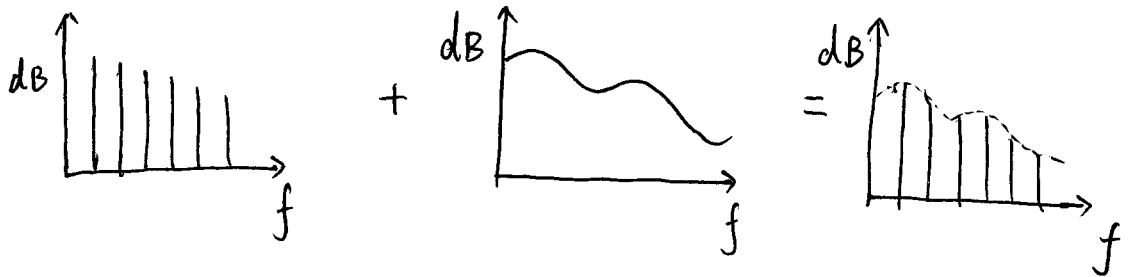i.e. the impulse response of the model has the same energy as the observed signal $x[n]$.

The error signal $e[n]$ is a our estimate of the excitation signal, where the resonance filter effect has been removed. Also called "whitened signal". Its spectrum has a flat envelope.

Cross synthesis: Combine $X_1[n]$'s excitation signal with $X_2[n]$'s filter to synthesize new signal $X_3[n]$.

① Calculate $H_1(z)$ of $X_1[n]$ using LPC.

② Calculate $e_1[n]$ by removing effect of $H_1(z)$

③ Calculate $H_2(z)$ of $X_2[n]$ using LPC

④ Filter $e_1[n]$ with $H_2(z)$ to synthesize $X_3[n]$.

Cepstrum ( cepstral domain representation of spectral envelope).

Basic idea:



$$e(t) \quad * \quad h(t) \quad = \quad x(t)$$
$$\Rightarrow \quad E(f) \quad \times \quad H(f) \quad = \quad X(f)$$
$$\Rightarrow \quad |E(f)| \quad \times \quad |H(f)| \quad = \quad |X(f)|$$
$$\Rightarrow \quad 20\log_{10}|E(f)| + 20\log_{10}|H(f)| = 20\log_{10}|X(f)|$$

View these spectra as ~~time do~~ "frequency-domain signal"

Excitation — periodic, high freq.

Filter — low freq.

Outcome — mixture

If we perform Fourier analysis on ~~X(t)~~ $20\log_{10}|X(f)|$ and transform it to a new domain, low coefficients will correspond to excitation, while high coefficients will correspond to the filter !
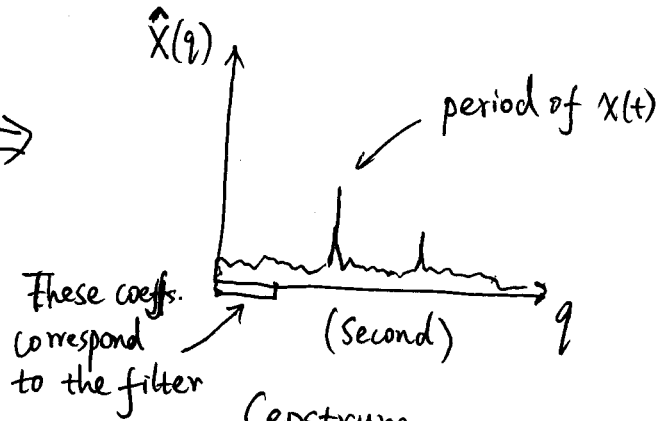
(Because Fourier transform is <u>linear</u> and it <u>separates low and high "freqs"</u>)

$\log_{10}|X(f)|$

(dB)

(Hz) $f$

$$\text{IFT} \Rightarrow$$

$\hat{X}(q)$

period of $x(t)$

These coeffs. correspond to the filter

(Second) $q$

spectrum

frequency

filtering

Cepstrum

quefrency

liftering

Digital implementation:

$$\hat{X}[q] = \text{IFFT}\left\{ \log | \text{FFT}[x[n]]| \right\}, \quad q = 0, \cdots, N-1.$$

It has the same length of spectrum and signal, N.

The first several (e.g., 20, 40) <u>cepstral coefficients</u> correspond to the resonance filter. The number is called the <u>order</u>.

In other words, we can reconstruct spectral envelope from thes cepstral coefficients by taking FFT.

Math:

$$\hat{X}[n] = \frac{1}{N} \sum_{k=0}^{N-1} a[k] e^{j2\pi kn/N} \qquad \text{where } a[k] = 20\log_{10}|X[k]|$$
$$\text{(symmetric)}$$

$$= \frac{1}{N} \sum_{k=0}^{N-1} a[k] \left\{ \cos\left(\frac{2\pi kn}{N}\right) + \underline{j \sin\left(\frac{2\pi kn}{N}\right)} \right\}$$

cancelled

$$= \frac{1}{N}\left( a[0] + (-1)^n a\left[\frac{N}{2}\right] \right) + \frac{2}{N} \sum_{k=1}^{N/2} a[k] \cos\left(\frac{2\pi kn}{N}\right)$$

↑ DC          ↑ Nyquist          ↑ Positive frequency.

$$= \text{DCT}\left\{ a[0:N/2] \right\} \qquad \text{(Discrete cosine transform of the positive frequency range log sp amplitude spectrum)}.$$

Another pespective to see cepstrum.

Log-amp spectrum can be approximated by a linear combination of several sinusoids, with coefficients of a cepstral coefficients.

$$a[k] \approx c_0 + \sqrt{2} \sum_{i=1}^{P-1} c_i \cos\left(2\pi i \frac{k}{N}\right) \quad , \text{ i.e.} \qquad (*)$$

$$\begin{bmatrix} a[0] \\ \vdots \\ a[N/2] \end{bmatrix} = \begin{bmatrix} 1 & \sqrt{2}\cos(2\pi 1 f_0) & \cdots & \sqrt{2}\cos(2\pi(p-1)f_0) \\ \vdots & \vdots & & \\ 1 & \sqrt{2}\cos(2\pi 1 f_{N/2}) & \cdots & \sqrt{2}\cos(2\pi(p-1)f_{N/2}) \end{bmatrix} \begin{bmatrix} c_0 \\ \vdots \\ \vdots \\ c_{p-1} \end{bmatrix}$$

where $f_k = k/N$.

$\downarrow$

M: First $p$ colums of a DCT matrix.
Colums are orthogonal.

Least-square solution:

$$\begin{bmatrix} c_0 \\ \vdots \\ c_{p-1} \end{bmatrix} = \underbrace{(M^T M)^{-1}}_{= 1/N} M^T \begin{bmatrix} a[0] \\ \vdots \\ a[N/2] \end{bmatrix} = \frac{1}{N} M^T \begin{bmatrix} a[0] \\ \vdots \\ a[N/2] \end{bmatrix}$$

$\uparrow$
$= 1/N$

∴ Cepstral coefficients are the least square solution to approximate log-amp spectrum using (*).

Mel-frequency Cepstral Coefficients (MFCC)

- broadly used in speech recognition, Speaker identification, timbre modeling.

- warp freq to log-freq in spectrum before taking IFFT.