

## Speech Signal Analysis

**DEF** (*Dithering*) Dithering adds noise to the signal which avoids the signal having zeros (which may be problematic when taking the log).

$$y[t] = x[t] + \epsilon[t] \quad \epsilon[t] \sim \mathcal{N}(0, \sigma^2)$$

**DEF** (*DC Offset*) Most processing techniques assume that the signal is centered around 0.

$$y[t] = x[t] - \frac{1}{T}\mu_x \quad \mu_x = \frac{1}{T} \sum_{t=1}^T x[t]$$

**DEF** (*Pre-emphasis*) Emphasize high frequency components.

$$y[t] = x[t] - \alpha x[t-1]$$

where  $\alpha$  is a constant (usually 0.97)

**DEF** (*Discrete Fourier Transform (DFT)*) The dot product of the signal with the sinusoids of the frequency (Fourier basis).

$$X[k] = \sum_{t=0}^{T-1} x[t] e^{i2\pi kt/T}$$

DFT decomposes the signal into frequency components. Commonly  $\mathcal{F}$  is used to denote the DFT operator.

**THM** (*Properties of DFT*)

- Linearity:  $\mathcal{F}(a_1x_1 + a_2x_2) = a_1\mathcal{F}(x_1) + a_2\mathcal{F}(x_2)$
- Shift: If  $y[t] = x[t-m]$ , then  $Y[k] = e^{i2\pi km/T} X[k]$

**DEF** (*Windowing*) Given a signal  $x[t]$ , a window function  $w[t]$  (Hamming, Hann, etc.)

$$y[t] = x[t]w[t]$$

**DEF** (*Short-Time Fourier Transform (STFT)*) Given a signal  $x[t]$ , a window function  $w[t]$ , and a hop size  $h$ , the STFT is given by:

$$X[k, m] = \sum_{t=0}^{T-1} x[t]w[t-m]e^{-i2\pi kt/T}$$

where  $m$  is the current frame index and  $k$  is the frequency index. This gives a complex spectrogram of the signal.

**DEF** (*Complex Spectrogram*) Given a frequency bin  $X[k] = a + bi$ ,

- Real:  $\mathcal{R}(X[k]) = a$
- Imaginary:  $\mathcal{I}(X[k]) = b$
- Magnitude:  $|X[k]| = \sqrt{a^2 + b^2}$
- Phase:  $\angle X[k] = \arccos \frac{a}{\sqrt{a^2 + b^2}}$
- Energy:  $P[k] = |X[k]|^2$

**DEF** (*Mel Spectrogram*) After obtaining the spectrogram,  $X[k]$ , and using the mel filter bank,  $H_n[k]$ , the mel spectrogram is given by:

$$Y[k, n] = \sum_{k=0}^{K-1} |X[k]|H_n[k]$$

**DEF** (*Mel Frequency Cepstral Coefficients (MFCCs)*) A common feature extraction technique for speech recognition which is able to capture speaker characteristics and phonetic content.

1. Extract the mel spectrogram  $Y[k, n]$
2. Apply the Discrete Cosine Transform (DCT) to the mel spectrogram.
3. Take the first  $p$  coefficients (usually 13).

## Hidden Markov Models

**DEF** (*Objective of ASR*) Given a sequence of acoustic feature vectors  $X$ , and  $W$  denotes a word sequence, ASR aims to find the most likely word sequence  $W^*$

$$W^* = \arg \max_W P(W|X)$$

**COR** (*Decomposition of  $P(W|X)$* ) We can decompose  $P(W|X)$  with Bayes' theorem:

$$P(W|X) = \frac{P(X|W)P(W)}{P(X)} \propto P(X|W)P(W)$$

$P(X|W)$  is the **Acoustic Model** and  $P(W)$  is the **Language Model**.

**DEF** (*Modelling the Acoustic Model with HMMs*) Commonly, the left-to-right HMM is used to model the acoustic model. As a word is composed of multiple phonemes, we can model each phoneme with a left-to-right HMM and then concatenate them to form a HMM for the word. For the phoneme HMMs, we typically use a left-to-right HMM with 3 states which as a consequence also enforces a minimum phone duration.

**DEF** (*The three fundamental problems of HMMs*)

1. **Likelihood** - Determine the overall likelihood of an observation sequence  $X = x_1, x_2, \dots, x_T$  given an HMM topology  $\mathcal{M}$
2. **Decoding** - Determine the most likely sequence of hidden states  $Q = q_1, q_2, \dots, q_T$  given an observation sequence  $X = x_1, x_2, \dots, x_T$  and an HMM topology  $\mathcal{M}$
3. **Training** - Given an observation sequence  $X = x_1, x_2, \dots, x_T$  and an HMM topology  $\mathcal{M}$ , determine the optimal state occupation probabilities.

**ALG** (*Forward Algorithm*) The **Likelihood** problem is solved by the **Forward Algorithm**. The forward probability  $\alpha_j(t)$  is the probability of observing the observation sequence  $x_1, \dots, x_t$  and being in state  $j$  at time  $t$ .

$$\alpha_j(t) = P(x_1, \dots, x_t, q_t = j | \mathcal{M})$$

This can be computed recursively:

- Initialisation:

$$\begin{aligned} \alpha_j(0) &= 1 & j &= 0 \\ \alpha_j(0) &= 0 & j &\neq 0 \end{aligned}$$

- Recursion:

$$\alpha_j(t) = \left( \sum_{i=0}^J \alpha_i(t-1)a_{ij} \right) b_j(x_t)$$

- Termination:

$$P(X|\mathcal{M}) = \alpha_E = \sum_{i=1}^J \alpha_i(T)a_{iE}$$