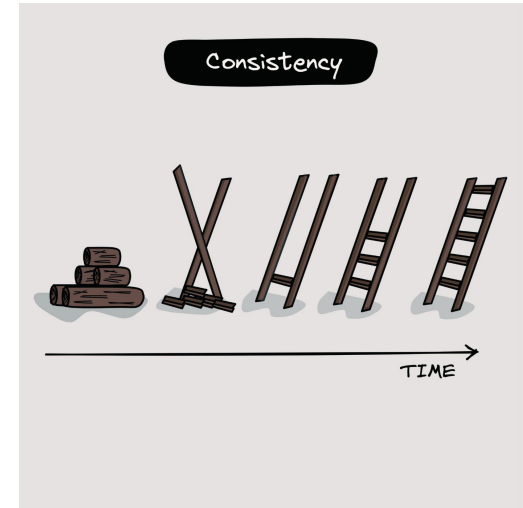


Case Study: Speech data and CNN

B.Tech. Data Science, NMIMS

By,

Bilal Hungund, Data Scientist, Halliburton



Audio Signal: (Automatic Speech Recognition)

Longitudinal vibration that produces vitality

Sound Wave:

Vibration signal produces by moving energy

Parameters

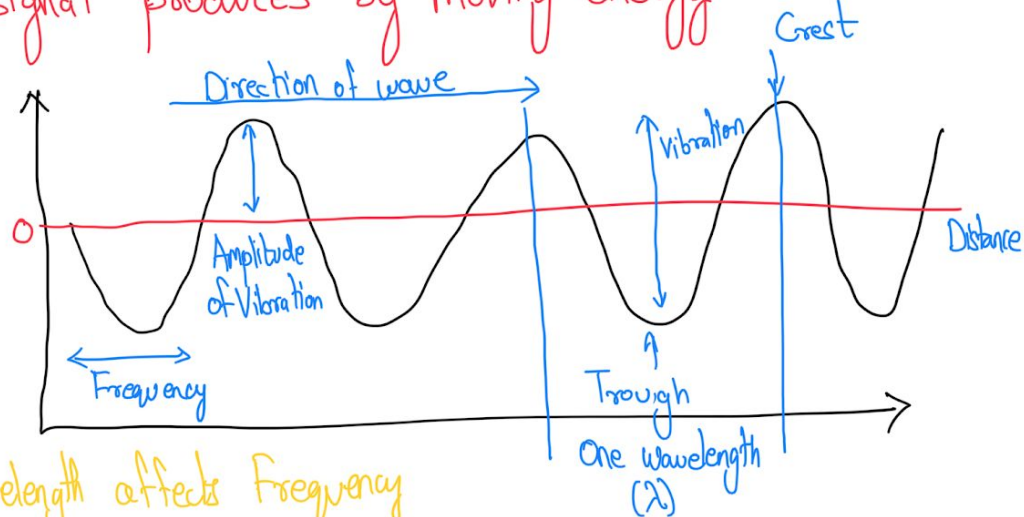
↳ Amplitude

Crest and Trough

Wavelength

Cycle

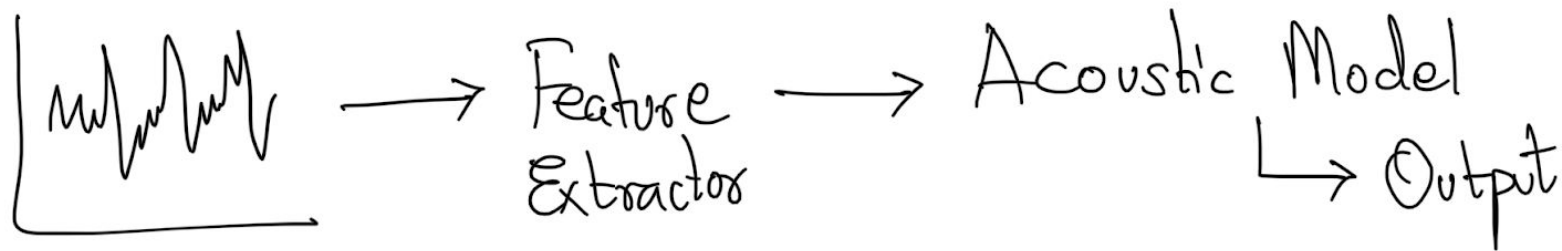
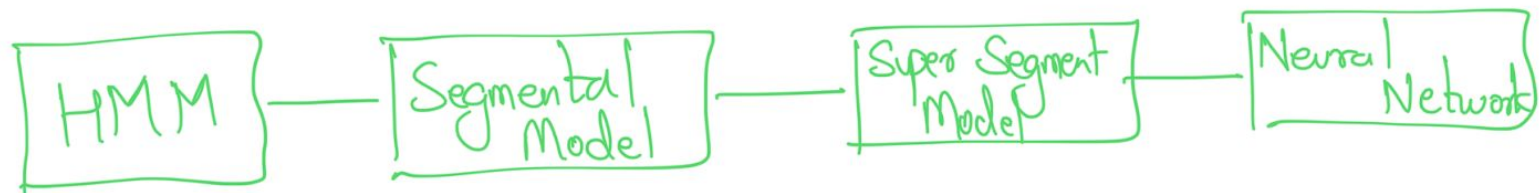
Frequency



Wavelength affects Frequency

Acoustic Modelling

→ Statistical Representation of computed feature vectors

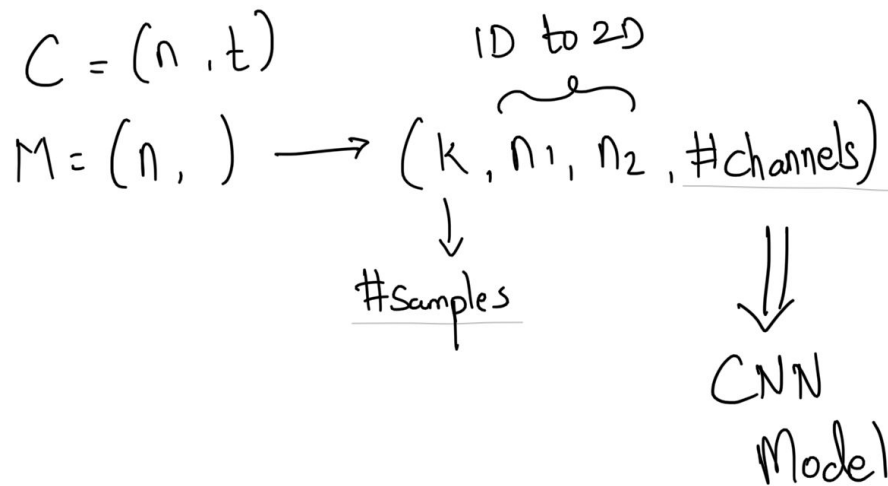
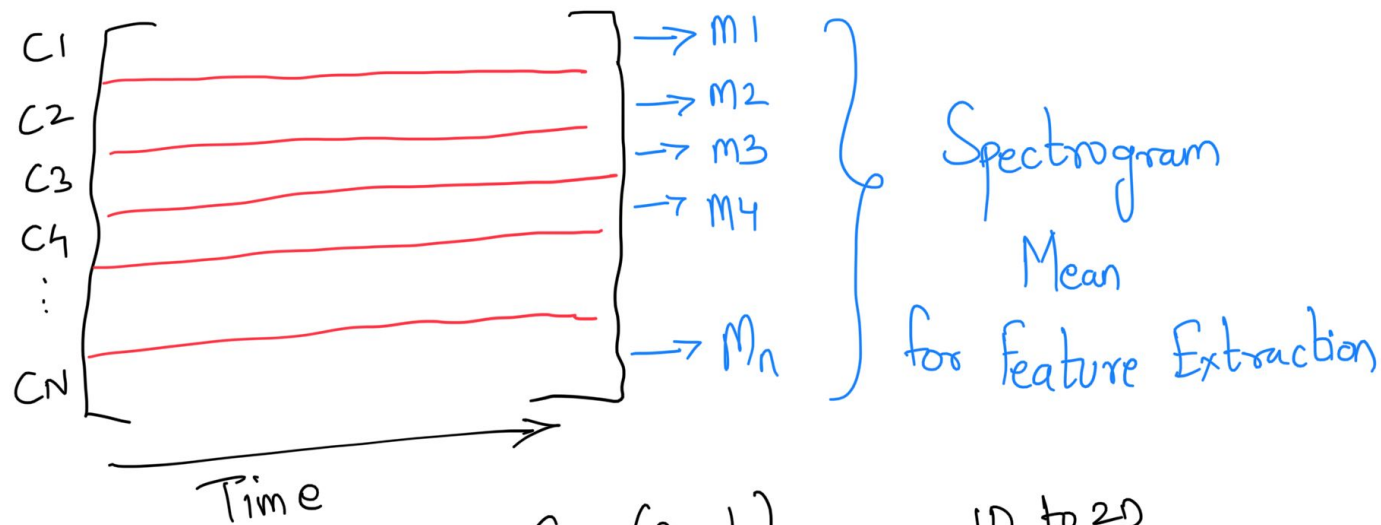


MFCC (Mel-frequency Cepstral Coefficients)

Mel Spectrogram

→ Spectrogram converted to Mel scale

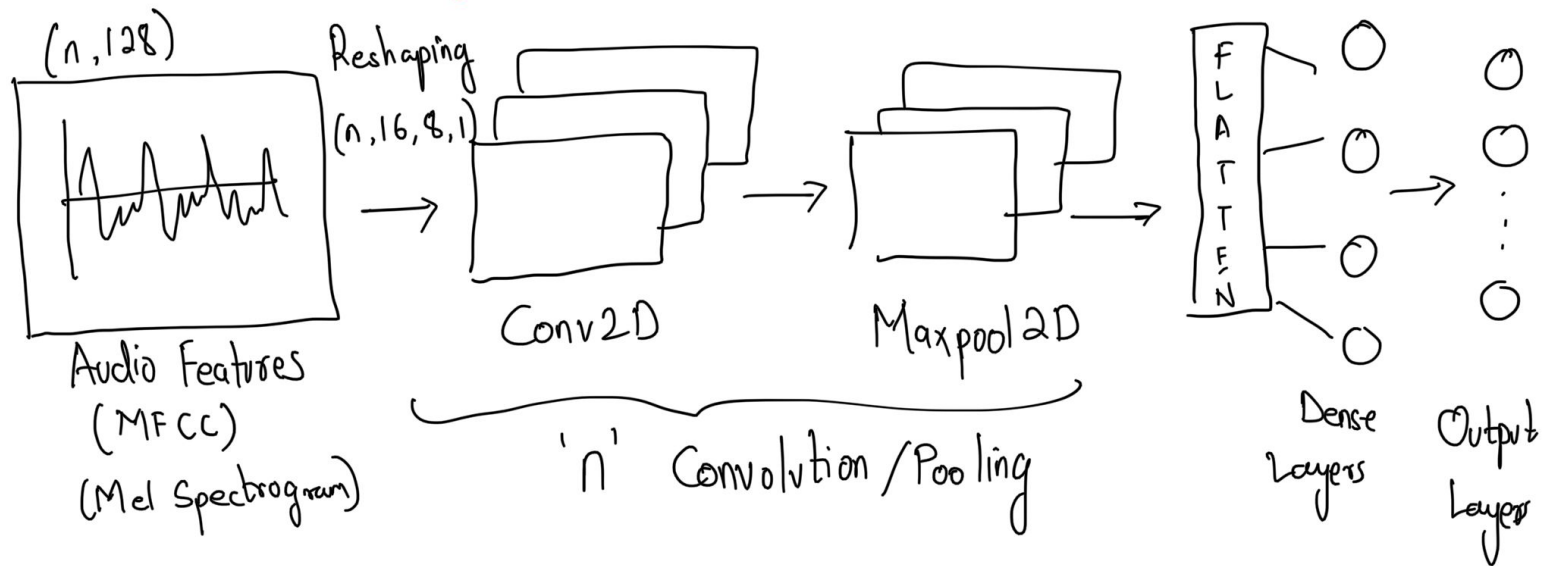
- Widely used in deep learning
- Powerful tool to extract the feature from speech
- Process includes: Fourier Transform, discrete cosine transforms and overlapping windows
- It helps for classification problems such as genre classification, disease detection related to speech and etc.



CNN in Speech Data

→ Create features using MFCCs & Mel Spectrogram

→ Average of matrix

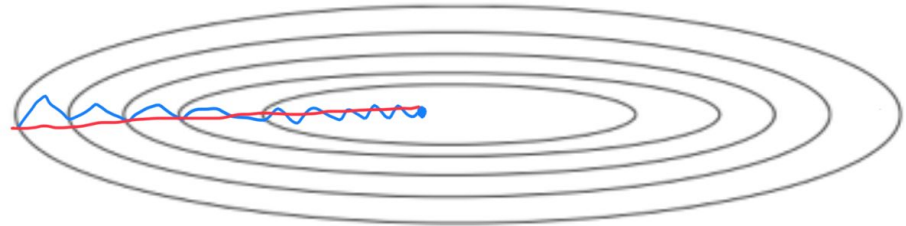
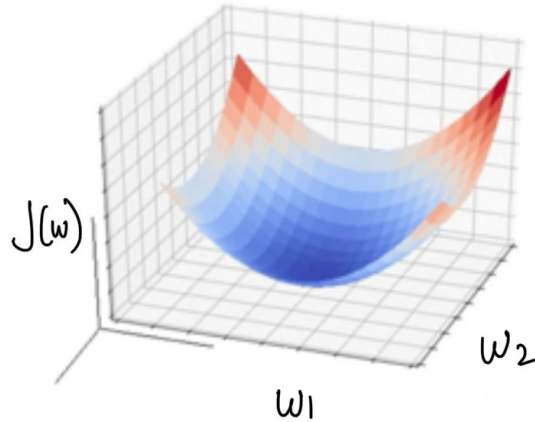


Adam Optimizer (Adaptive moment estimation)

↳ Momentum
↳ RMSprop

Exponentially
weighted moving
average

Momentum:



Momentum:

$$V_{dw} = \beta V_{dw} + (1-\beta) dW$$

$$V_{dB} = \beta V_{dB} + (1-\beta) dB$$

$$W := W - \eta V_{dw}$$

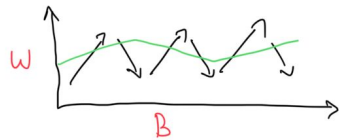
$$B := B - \eta V_{dB} \quad [\eta = \text{learning rate}]$$

$$W := W - \eta dW$$

$$B := B - \eta dB$$

$$\beta = 0.9$$

RMS prop (Root mean squared proportion)



$$S_{dw} = \beta S_{dw} + (1-\beta) (dW)^2$$

$$S_{dB} = \beta S_{dB} + (1-\beta) (dB)^2$$

$$W := W - \eta \frac{dW}{\sqrt{S_{dw} + \epsilon}}$$

$$\beta = 0.999$$

$$\epsilon = 10^{-8}$$

$$B := B - \eta \frac{dB}{\sqrt{S_{dB} + \epsilon}}$$

Adam:

$$W := W - \eta \frac{V_{dw}}{\sqrt{S_{dw} + \epsilon}}$$

$$B := B - \eta \frac{V_{dB}}{\sqrt{S_{dB} + \epsilon}}$$

Momentum β will be β_1 | RMS prop β will be β_2