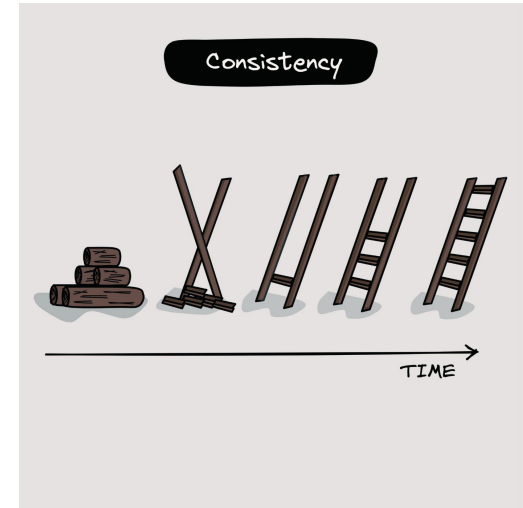# Case Study: Speech data and CNN

M.Tech. Data Science, Second Year, NMIMS

By,

Bilal Hungund, Data Scientist, Halliburton

# Convolution & operation



Filter (Weights)

$\times$

3×3 Filter

6×6 image

$n - f + 1$

$= 6 - 3 + 1$

$= 4$

Output Size

Output Without Padding

Stride = 1 (step size)

4×4

## With Padding

Padding

filter

4×4 $\longrightarrow$ 6×6

Padding

4×4

# Pooling

| 25 | 48 |
|----|----|
| 192 | 10 |

| 192 | 110 |
|-----|-----|
| 50 | 47 |

| 25 | 48 | 11 | 58 |
|----|----|----|----|
| 192 | 10 | 20 | 110 |
| 38 | 0 | 9 | 31 |
| 50 | 8 | 23 | 47 |

| 11 | 58 |
|----|----|
| 20 | 110 |

Pooling $\Rightarrow$

| 38 | 0 |
|----|----|
| 50 | 8 |

| 69 | 50 |
|----|----|
| 22 | 28 |

Average Pooling

Stride = 2
(Recommended for Pooling)

| 9 | 31 |
|---|----|
| 23 | 47 |

# Convolution Neural Network (CNN)
## for Classification

Conv-n
Conv-2
Conv-1

Pooling

D
N
N

$(i/2, j/2)$

$(i,j)$

Filters: n

tf.keras.layers.Conv2D
tf.keras.activations.*
tf.keras.layers.MaxPool2D

1) Convolution: Filters to generate feature maps

2) Non-linearity: often relu

3) Backpropagation    4) Pooling: Downsampling feature maps

# Audio Signal: (Automatic Speech Recognition)

Longitudinal vibration that produces vitality

# Sound Wave:

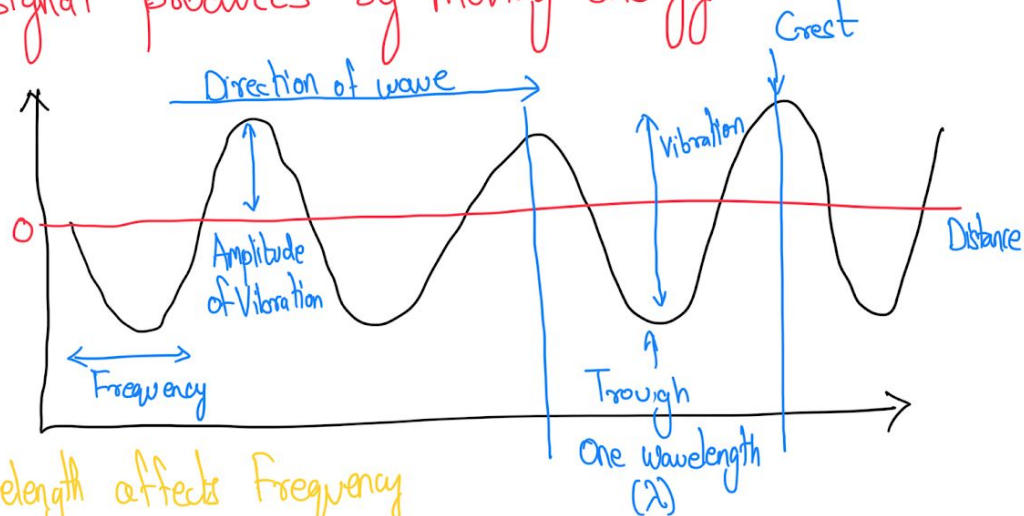Vibration signal produces by moving energy

Parameters
↳ Amplitude
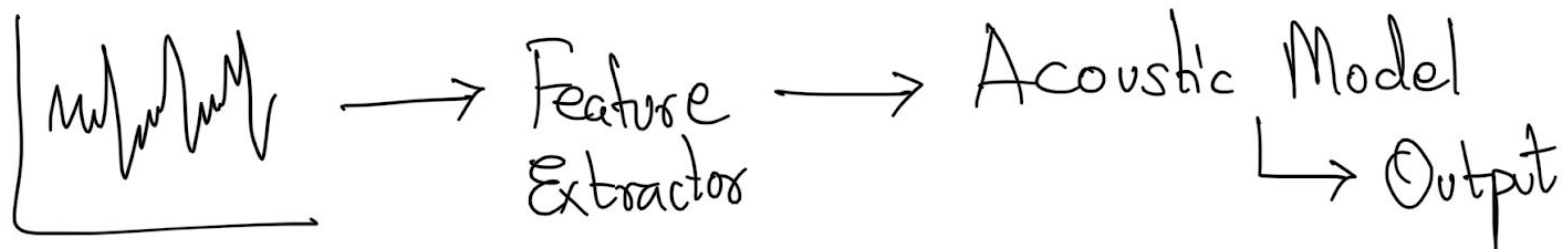Crest and Trough
Wavelength
Cycle
Frequency

Direction of wave →

0

Crest

Vibration

Amplitude of Vibration

Frequency

Distance

Trough

One Wavelength (λ)

Wavelength affects Frequency

# Acoustic Modelling

→ Statistical Representation of computed feature vector

```
┌─────────┐     ┌──────────┐     ┌──────────────┐     ┌──────────┐
│  HMM    │─────│ Segmental│─────│ Super Segment│─────│ Neural   │
│         │     │  Model   │     │    Model     │     │ Network  │
└─────────┘     └──────────┘     └──────────────┘     └──────────┘
```



Feature Extractor → Acoustic Model ↳ Output

# MFCC (Mel-frequency Ceptral Coefficients)

Mel Spectrogram

→ Spectrogram Converted to Mel scale

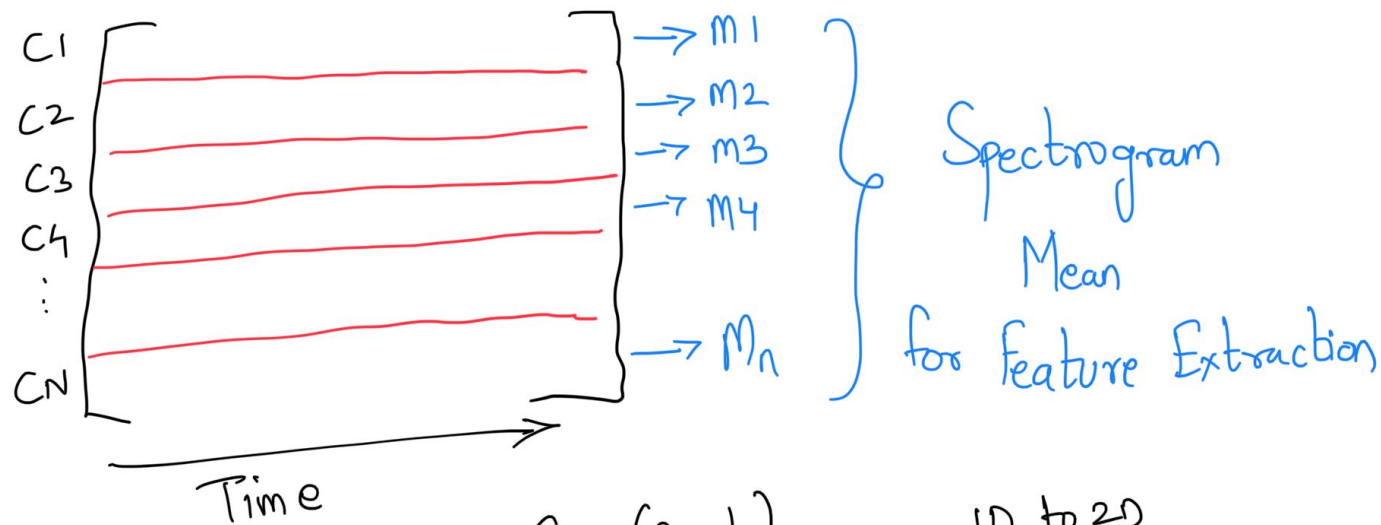→ Widely used in deep learning

→ 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.

C1

C2

C3

C4

...

CN

Time

$\rightarrow$ M1

$\rightarrow$ M2

$\rightarrow$ M3

$\rightarrow$ M4

$\rightarrow$ $M_n$

Spectrogram

Mean

for Feature Extraction

$C = (n, t)$

$M = (n, \ ) \longrightarrow (k, \overbrace{n_1, n_2}^{\text{1D to 2D}}, \underline{\#channels})$

$\downarrow$

$\underline{\#Samples}$

$\Downarrow$

CNN Model

# CNN in Speech Data

→ Create features using MFCCs & Mel Spectrogram

→ Average of matrix



(n, 128)

Reshaping

(n, 16, 8, 1)

Audio Features
(MFCC)
(Mel Spectrogram)

Conv2D

Maxpool2D

FLATTEN

'n' Convolution/Pooling

Dense Layers

Output Layers