

Week - 6 - DLP

↳ Speaker diarization - Intro

- Speaker diarisation → a applicaⁿ in speech technology.
multiple speakers are there. $s(t)$ (what was said)
→ t text . (ASR) (audio to text)
also called as transcription.

- instead of 1 big text str.

$\left. \begin{matrix} s_1. \\ s_2. \\ s_3. \end{matrix} \right\}$ like this would be better

Speaker diarisaⁿ (who spoke when)

2. ↳ we will use speaker identification model to achieve this.

(whisper^{ASR} model) ← by Open AI.
dominated by American English.

Eg.  → ASR 2.03 - 2.08 —

2.09 - 2.17 —



Speaker Identification model

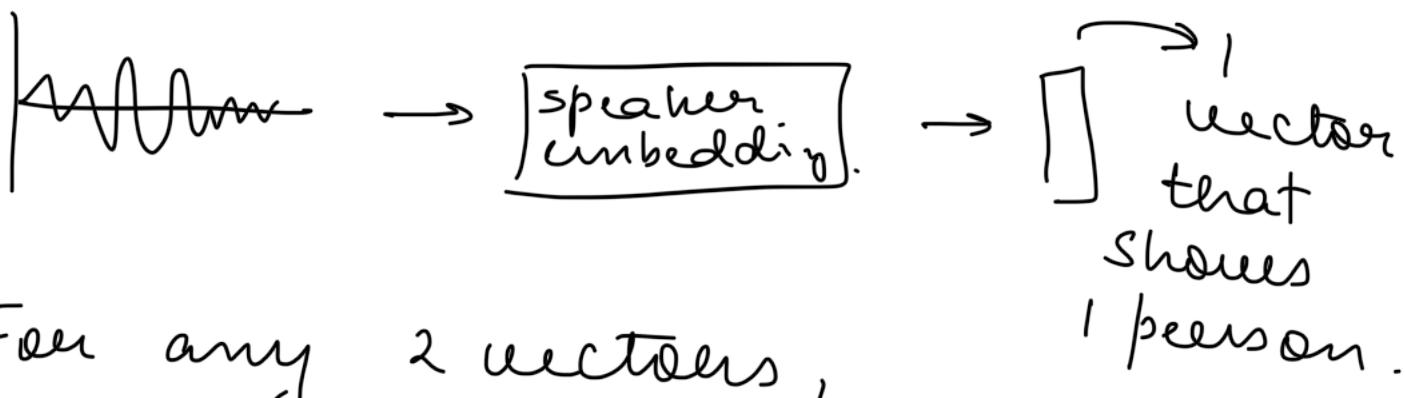
↳ Speaker embedding extender

ECAPA-TDNN

Model will
be used.

Speech
Toolkit
Kaldi
Espresso
Speech
Brain

L2 Speaker Identification



For any 2 vectors,

[] & [] we can calculate some kind of distance (euclidean dist.) to show the degree of closeness.

Or

we can use the angle b/w 2 vectors to explain closeness.

$$\cos \theta = \frac{x^T y}{\|x\| \cdot \|y\|}$$

$\|x\| \cdot \|y\|$ → cosine distance



calculate the cosine dist. b/w these embeddings.

if, $\langle s_1, s_2 \rangle < \langle s_1, s_3 \rangle$

then, s_1 is closer to s_2 than s_3 .

- for our ease, we will assume that there are only n no. of speakers in my recording.

L3 Clustering Techniques

- Agglomerative Clustering \rightarrow keep merging clusters, until I hit n (desired no. of clusters).

Speaker \rightarrow I know there are people only.
Initially, we have

7 segments \longrightarrow 7 embeddings
 E_1, \dots, E_7

calculate,

$$\begin{array}{c} \langle E_1, E_2 \rangle \\ \langle E_2, E_3 \rangle \end{array} \dots \dots \begin{array}{c} \langle E_1, E_7 \rangle \\ \langle E_2, E_7 \rangle \end{array}$$

$$\langle E_6, E_7 \rangle$$

find the minimum cosine distance

e.g. $\langle E_4, E_7 \rangle$ was max.

$$\hookrightarrow \langle E'_4, E'_4 \rangle$$

and repeat this. \rightarrow finally we reach a stage where there are only 2 embeddings. \therefore , 2 speakers.

L4 demo

speechbrain == 0.5.16
faster-whisper
pyannote-audio
whisper

} use lib.

- Convert any format to .wvav using ffmpeg.