# Improving Unsupervised Acoustic Word Embeddings using Speaker and Gender Information

Lisa van Staden, Herman Kamper

31 January 2020

UNIVERSITEIT
iYUNIVESITHI
STELLENBOSCH
UNIVERSITY

100
1918·2018

Popular methods for speech processing rely on transcribed speech.



i   had   to   think      of      some   example speech

since      speech  recognition  is   really cool

Popular methods for speech processing rely on transcribed speech.



i    had   to   think      of       some   example speech

since      speech  recognition  is   really  cool

Obtaining transcriptions is expensive and not always possible.

We don't always need to predict text labels:

We don't always need to predict text labels:

- Query-by-Example Search: search speech using speech.

We don't always need to predict text labels:

- Query-by-Example Search: search speech using speech.

- Unsupervised Term Discovery: Discover repeating patterns in speech.
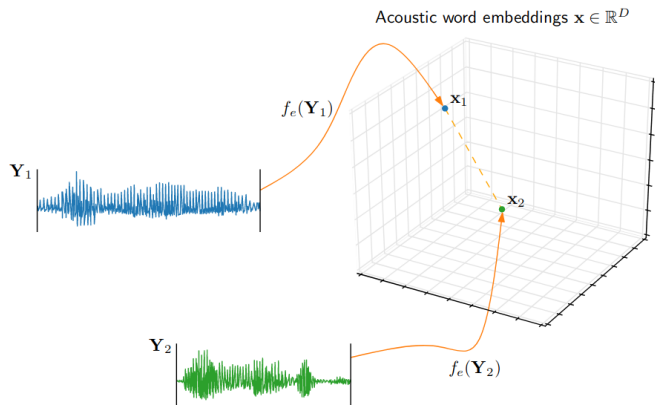
These tasks require comparing speech segments.

The conventional method is Dynamic Time Warping.

These tasks require comparing speech segments.

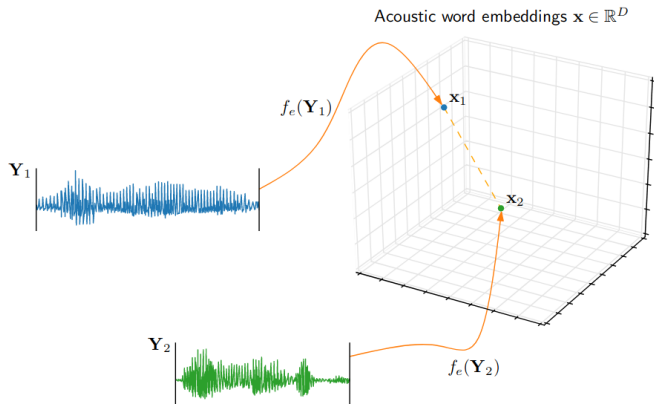The conventional method is Dynamic Time Warping.

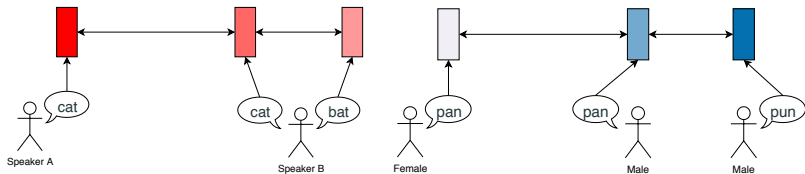- Computationally expensive.

Acoustic word embeddings $\mathbf{x} \in \mathbb{R}^D$

$f_e(\mathbf{Y}_1)$

$\mathbf{x}_1$

$\mathbf{Y}_1$

$\mathbf{x}_2$

$\mathbf{Y}_2$

$f_e(\mathbf{Y}_2)$

Acoustic word embeddings $\mathbf{x} \in \mathbb{R}^D$

$f_e(\mathbf{Y}_1)$

$\mathbf{x}_1$

$\mathbf{Y}_1$

$\mathbf{x}_2$

$\mathbf{Y}_2$

$f_e(\mathbf{Y}_2)$

We want to map speech to these representation without using labels.

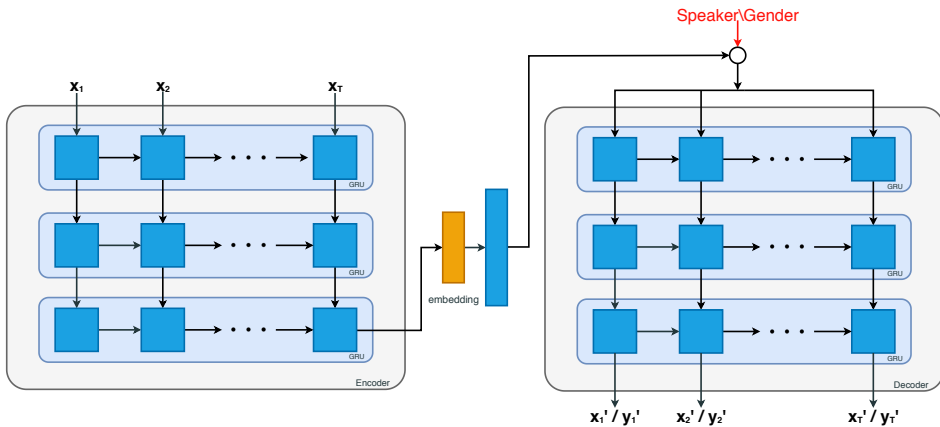Acoustic properties of speech from different speakers/genders differ.
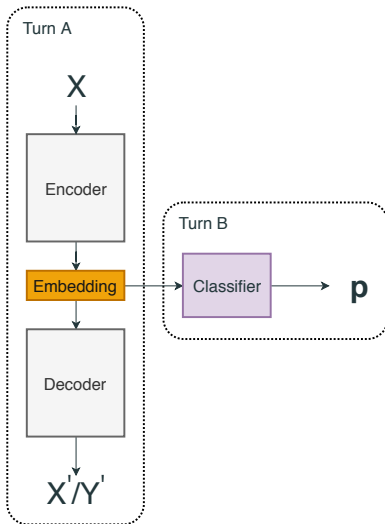


We want embeddings to be robust.

Turn A

X

Encoder

Turn B

Embedding → Classifier → p

Decoder

X'/Y'

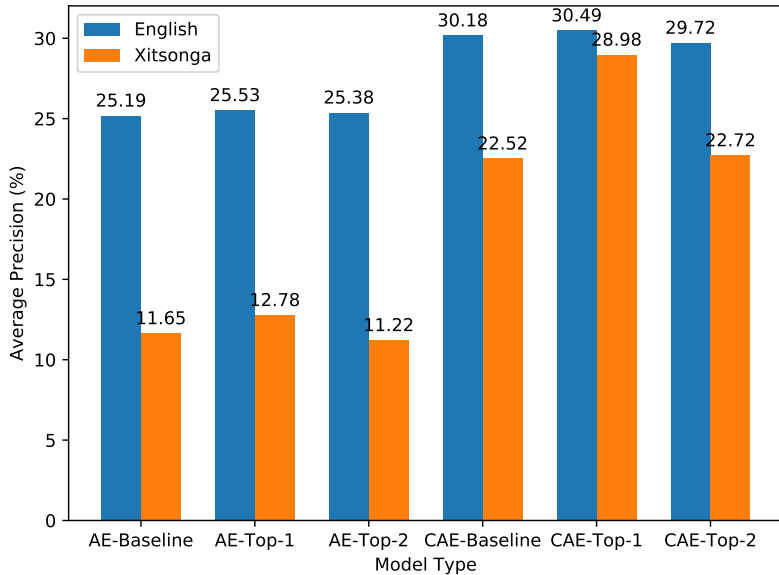z → Linear ReLU Linear ReLU Dropout Linear Softmax → p

■ Linear ■ ReLU ■ Dropout ■ Softmax

Use the same-different task to evaluate AWEs:

- Measure if AWEs are similar given a threshold.

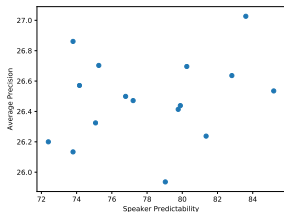- Calculate area under Precision vs Recall curve.

# Results

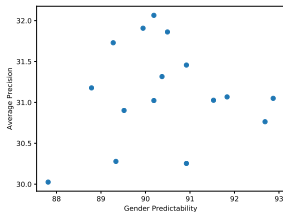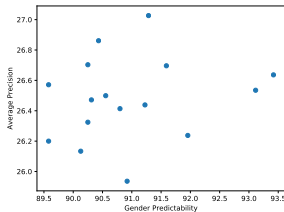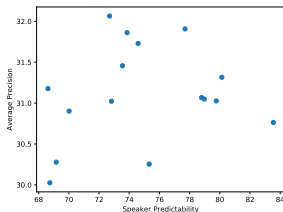Analyse if the speaker and gender information has decreased:

- Use speaker/gender classifier model.

- Evaluate accuracy.

# Average Precision vs Speaker/Gender Predictability

- English data shows marginal improvement by incorporating speaker information.

- English data shows marginal improvement by incorporating speaker information.

- Best Xitsonga model shows 22% improvement.

## Conclusions

- English data shows marginal improvement by incorporating speaker information.

- Best Xitsonga model shows 22% improvement.

- It's difficult to remove speaker and gender information.

## Conclusions

- English data shows marginal improvement by incorporating speaker information.

- Best Xitsonga model shows 22% improvement.

- It's difficult to remove speaker and gender information.

- Future work …