

# Triplet loss functions in speaker recognition systems

---

Francisco Javier Sáez Maldonado

May 10, 2022

Máster en Ciencia de Datos

*Escuela Politécnica Superior  
Universidad Autónoma de Madrid*

- **Task:** Speaker Recognition
  - 'Closed-set' vs 'Open-set'

- **Task:** Speaker Recognition
  - 'Closed-set' vs 'Open-set'
- **Data:** VoxCeleb challenge dataset.

## In defence of metric learning for speaker recognition

*Joon Son Chung, Jaesung Huh, Seongkyu Mun, Minjae Lee, Hee Soo Heo,  
Soyeon Choe, Chiheon Ham, Sunghwan Jung, Bong-Jin Lee, Icksang Han*

Naver Corporation, South Korea

joonson.chung@navercorp.com

### Abstract

The objective of this paper is ‘open-set’ speaker recognition of unseen speakers, where ideal embeddings should be able to condense information into a compact utterance-level representation that has small intra-speaker and large inter-speaker distance.

A popular belief in speaker recognition is that networks trained with classification objectives outperform metric learning methods. In this paper, we present an extensive evaluation of most popular loss functions for speaker recognition on the VoxCeleb dataset. We demonstrate that the vanilla triplet loss shows competitive performance compared to classification-based losses, and those trained with our proposed metric learning objective outperform state-of-the-art methods.

popular due to their ease of implementation and good performance [17, 18, 19, 20, 21, 22, 23, 24]. However, training with AM-Softmax and AAM-Softmax has proven to be challenging since they are sensitive to the value of scale and margin in the loss function.

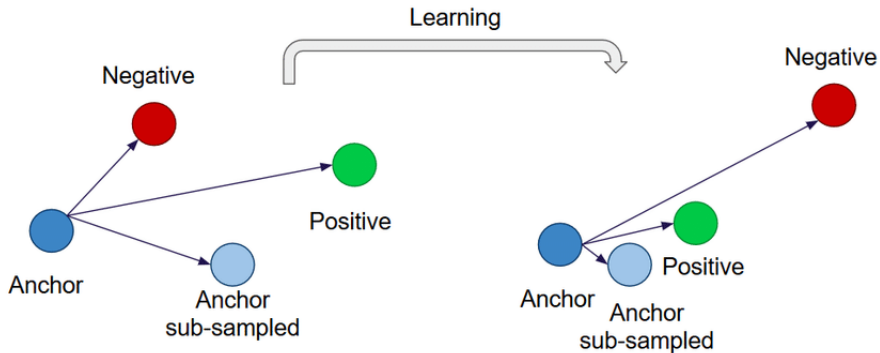
Metric learning objectives present strong alternatives to the prevailing classification-based methods, by learning embeddings directly. Since open-set speaker recognition is essentially a metric learning problem, the key is to learn features that have small intra-class and large inter-class distance. Contrastive loss [25] and triplet loss [26] have been demonstrated promising performance on speaker recognition [27, 28] by optimising the distance metrics directly, but these methods require careful pair or triplet selection which can be time consuming and performance sensitive.

## Softmax is not enough

- **Softmax:** 
$$L_S = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^C e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}}$$

## Softmax is not enough

- **Softmax:**  $L_S = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{w}_{y_i}^T \mathbf{x}_i + b_{y_i}}}{\sum_{j=1}^C e^{\mathbf{w}_j^T \mathbf{x}_i + b_j}}$



## Tools

---

# Conclusions

- Usage of two general machine learning techniques applied to natural language processing.
- Applying this method has a similar effect to regularization.
- Results are promising in Language modeling, but not so much in Neural Machine Translation.



Thank you for your **attention**.

