



(b) Open-set verification

Open-set Speaker Verification:

(determines whether a pair of utterances belongs to the same person)

- Speaker identities in testing set are usually disjoint from the ones in training set, which makes the speaker verification more challenging yet closer to practice.
- Since it is impossible to classify testing utterances to known identities in training set, we need to map speakers to a discriminative feature space.
- > In this scenario, open-set speaker verification is essentially a metric learning problem, where the key is to learn discriminative large-margin features.

Traditional Length Normalization

Length normalization on i-vector has been the de facto standard before back-end modeling.

For open-set SV task, cosine similarity or length normalization followed by probabilistic linear discriminant analysis (PLDA) scoring modeling is widely used to get the final pairwise scores.

The cosine similarity is a similarity measure which is independent of magnitude, it can be seen as the length-normalized version of inner-product of two vectors..





Once speaker embeddings (such as x-vectors) are extracted, just the same as in i-vector approach, cosine similarity or length normalization followed by PLDA is commonly adopted to get the final pairwise scores.

Length normalization in typical speaker embedding approach

peaker Verification System

 Li^1 an. China Guangzhou, China

Deep Length Normalization

Irn the deep speaker embeddings being length-normalized in an enu-to-enu manner within common classification network? We add a length normalization layer followed by a scale layer before the $\mathbf{v}_i = \alpha \times \frac{f(\mathbf{x}_i)}{1 - 1}$ outnut laver of the common clas []]-[] feature sequence Deep length normalization Con^{*} layer Speaker embeddi eature sequence **Experimental Results and Discussion** Table 1: Baseline end-to-end system architecture > The model is trained with a mini-batch size of 128, using typical stochastic gradient descent with momentum 0.9 and weight decay 1e-4. > The learning rate is set to 0.1, 0.01, 0.001 and is switched when the training loss Layer Channels Blocks Output size Downsample plateaus. Conv1 $64 \times L$ False 16 > For each training step, an integer L within [300,800] interval is randomly Res1 16 $64 \times L$ False generated, and each data in the mini-batch is cropped or extended to L frames. Res2 $32 \times$ True 32 > After model training finished, the 128-dimensional speaker embeddings are 64 Res3 $16 \times$ True extracted after the penultimate layer of neural network. True Res4 $8 \times -$ 128 Score 128 Average pool -128 FC (embedding) -Table 3: Verification performance on VoxCeleb1 for various speaker categories Output scale parameter α (lower is better) minDCF10⁻³ EER(%)System Description minDCF10⁻² **Experiments on different** α 0.553 0.713 5.48 Deep embedding baseline • We first investigate the setting of scale parameter α . For those systems in Table 0.922 10.18 3 and Fig. 4, the cosine similarity or equivalently L2-normalized inner-product is fixed $\alpha = 1$ 0.967 0.601 0.828 6.36 adopted to measure the similarities between speaker embeddings. fixed $\alpha = 4$ • From Fig. 4, we can observe the proposed L2-normalized deep embedding 0.515 0.687 5.49 fixed $\alpha = 8$ 0.586 5.01 0.475 fixed $\alpha = 12$ outperforms the baseline system significantly. fixed $\alpha = 16$ 0.499 0.596 5.32 • The performance is poor when α is too small and stable with α is higher. The 0.637 0.503 5.46 fixed $\alpha = 20$ 0.502 0.638 5.54 best α in our experiment is 12. fixed $\alpha = 24$ 5.52 fixed $\alpha = 28$ 0.497 0.640 0.486 0.599 5.60 trained $\alpha = 26.1$ Table 2: Voxceleb1 open-set verification task performance, in comparing the effect of our introduced deep length normalization strategy and traditional extra length normalization step (lower is better)

System Description	Deep L_2 -norm	Traditional L ₂ -norm	Similarity Metric	minDCF10 ⁻²	minDCF10 ⁻³	EER(%
i-vector + inner-product	N/A	X	inner-product	0.736	0.800	13.80
i-vector + cosine	N/A	\checkmark	inner-product	0.681	0.771	13.80
i-vector + PLDA	N/A	X	PLDA	0.488	0.639	5.48
i-vector + L_2 -norm + PLDA	N/A	\checkmark	PLDA	0.484	0.627	5.41
Deep embedding + inner-product	X	X	inner-product	0.758	0.888	7.42
Deep embedding+ cosine	×	\checkmark	inner-product	0.553	0.713	5.48
Deep embedding+ PLDA	×	X	PLDA	0.524	0.739	5.90
Deep embedding + L_2 -norm + PLDA	×	\checkmark	PLDA	0.545	0.733	5.21
<i>L</i> ₂ -normalized deep embedding + inner-product	1	X	inner-product	0.475	0.586	5.01
L_2 -normalized deep embedding + PLDA	✓	X	PLDA	0.525	0.694	4.74

- system achieves the best minDCF of 0.475, 0.586 and EER of 5.01%, which

We further compare the effect of deep length normalization strategy and traditional extra length normalization in the whole SV pipeline. The results are shown in Table 2.

- achieves the best performance.
- When it turns into L2-normalized deep speaker embedding more extra length normalization step.
- function introduced by PLDA.







• No matter in i-vector or baseline deep speaker embedding systems, extra length normalization step followed by PLDA scoring

systems, since the speaker embeddings extracted from neural network have already been normalized to unit length, we need no

 In testing stage, a simple inner-product achieves the best performance, even slightly better than the PLDA scoring result. It might be the reason that our L2-normalized speaker embedding is highly optimized, which could incompatible with the objective