# RESEARCH

# A Simulation Study on Optimal Scores for Speaker Recognition

Dong Wang

Correspondence: wangdong99@mails.tsinghua.edu.cn Center for Speech and Language Technologies, Tsinghua University, Beijing, China Full list of author information is available at the end of the article

# Abstract

In this article, we conduct a comprehensive simulation study for the optimal scores of speaker recognition systems that are based on speaker embedding. For that purpose, we first revisit the optimal scores for the speaker identification (SI) task and the speaker verification (SV) task in the sense of minimum Bayes risk (MBR), and show that the optimal scores for the two tasks can be formulated as a single form of normalized likelihood (NL). We show that when the underlying model is linear Gaussian, the NL score is mathematically equivalent to the PLDA likelihood ratio (LR), and the empirical scores based on cosine distance and Euclidean distance can be seen as approximations of this linear Gaussian NL score under some conditions.

Based on the unified NL score, we conducted a comprehensive simulation study to investigate the behavior of the scoring component on both the SI task and SV task, in the case where the distribution of the speaker vectors perfectly matches the assumption of the NL model, as well as the case where some mismatch is involved. Importantly, our simulation is based on the statistics of speaker vectors derived from a practical speaker recognition system, hence reflecting the behavior of the NL scoring in real-life scenarios that are full of imperfection, including non-Gaussianality, non-homogeneity, and domain/condition mismatch.

Keywords: Speaker recognition; Speaker embedding; Scoring

# 1 Introduction

With decades of investigation, speaker recognition has achieved significant performance, and has been deployed in a wide range of practical applications [10, 26, 48]. Speaker recognition research concerns two tasks: speaker identification (SI) that identify the true speaker from a set of candidates, and speaker verification (SV) that tests if an alleged speaker is the true speaker. The performance of SI systems is evaluated by identification rate (IDR), the percentage of the trials whose speakers are correctly identified. SV systems require a threshold to decide whether accepting the speaker or not and the performance is evaluated by equal error rate (EER), to represent the trade-off between fail to accept and fail to reject.

Modern speaker recognition methods are based on the concept of *speaker embedding*, i.e., representing speakers by fixed-length continuous *speaker vectors*. This embedding is traditionally based on statistical models, in particular the i-vector model [17]. Recently, deep learning methods gained much attention and embedding based on deep neural nets (DNN) becomes popular [35, 60]. With the efforts from multiple research groups, deep speaker embedding models have been significantly improved by comprehensive architectures [14, 29], smart pooling approaches [8, 11, 44, 66], task-oriented objectives [3, 18, 22, 34, 62, 70], and carefully designed training schemes [39, 57, 64]. As a result, the deep embedding approach has achieved state-of-the-art performance [51]. Among various deep embedding architectures, the x-vector model is the most popular [55].

A key component of the speaker embedding approach is how to score a trial. Numerous empirical evidence has shown that the likelihood ratio (LR) derived by probabilistic linear discriminant analysis (PLDA) [28, 46] works well in most situations, and when the computational resource is limited, the cosine distance is a reasonable substitution. In some circumstances in particular on SI tasks, the Euclidean distance can be used. In this article, we revisit the scoring methods for speaker recognition from the perspective of minimum Bayes risk (MBR). The analysis shows that for both the SI and SV tasks, the MBR optimal score can be formulated as a single form  $\frac{p_k(\boldsymbol{x})}{p(\boldsymbol{x})}$ , which we call a normalized likelihood (NL) score. In the NL score,  $p_k(\boldsymbol{x})$ is the likelihood term that represents the probability that the test utterance  $\boldsymbol{x}$  belongs to the target class k, and  $p(\mathbf{x})$  is a normalization term that represents the probability that  $\boldsymbol{x}$  belongs to all possible classes. We will show that the NL score is equivalent to PLDA LR, in the case where the speaker vectors are modeled by a linear Gaussian and the target class is represented by finite enrollment utterances. We will also show that under some conditions, the empirical scores based on cosine distance and Euclidean distance can be derived from the linear Gaussian NL score.

Based on the unified formulation of the NL score, we conducted a comprehensive simulation study on the performance bound of a speaker recognition system, on both the SI and SV tasks. In particular, by imitating the statistical properties of speaker vectors derived from a *real* recognition system, our simulation gained deep understanding of a modern speaker recognition system, for instance the upper bound of its performance, and its behavior with real-life imperfection, including non-Gaussianality, non-homogeneity, training-deployment domain mismatch and enrollment-test condition mismatch. To the best knowledge of the author, this is the first comprehensive simulation study on the scoring component of modern speaker recognition systems. Note that the NL formulation is a prerequisite for the simulation study: it not only allows using the same score to investigate the behavior of both the SI and SV systems, but also offers the possibility to decompose the scoring model into separate components (by using different statistical models), which is important when we analyze the domain and condition mismatch.

It should be noted that the NL formulation is not new and may trace back to the LR scoring method with the Gaussian mixture model-Universal background model (GMM-UBM) framework [49]. Within the speaker embedding framework, the NL form was derived by McCree et. al. [6, 42] from the hypothesis test view (the one used for PLDA inference). Our derivation is based on the MBR decision theory, which directly affirms the optimum of the NL score.

The rest of the paper is organized as follows: Section 2 will revisit the MBR optimal scoring theory and propose the NL score. Section 3 presents the simulation results. Some discussions are presented in Section 4 and the entire paper is concluded in Section 5.

# 2 Theory and methods

2.1 MBR optimal decision and normalized likelihood

It is well known that an optimal decision for a classification task should minimize the Bayes risk (MBR):

$$k^* = \arg\min_k \sum_j \ell_{jk} p(j|\boldsymbol{x}) \tag{1}$$

where  $\boldsymbol{x}$  is the observation,  $\ell_{jk}$  is the risk taken when classifying an observation from class j to class k. In the case where  $\ell_{jk}$  is 0 for j = k and a constant c for any  $j \neq k$ , the MBR decision is equal to selecting the class with the largest posterior probability:

$$k^* = \arg\max_k p(k|\boldsymbol{x}). \tag{2}$$

We call this result the *MAP principle*. We will employ this principle to derive the optimal score for the SI and SV tasks in speaker recognition.

# 2.1.1 MBR optimal score for SI

In the SI task, our goal is to test K outcomes  $\{H_k: \boldsymbol{x} \text{ belongs to class } k\}$  and make the decision which outcome is the most probable. Following the MAP principle, the MBR optimal decision is to choose the k-th outcome that obtains the maximum posterior:

$$k^* = \arg\max_k p(H_k | \boldsymbol{x}) = \arg\max_k p_k(\boldsymbol{x}) p(k),$$
(3)

where k indexes the classes, and  $p_k(\mathbf{x})$  represents the likelihood of  $\mathbf{x}$  in class k. In most cases, there is no preference for any particular class and so the prior p(k) for each class k shall be equal. We therefore have:

$$k^* = \arg\max_k p_k(\boldsymbol{x}). \tag{4}$$

It indicates that MBR optimal decisions can be conducted based on the likelihood  $p_k(\boldsymbol{x})$ . In other words, the likelihood is MBR optimal for the SI task.

## 2.1.2 MBR optimal score for SV

For the SV task, our goal is to test two outcomes and check which one is more probable: {  $H_0$ :  $\boldsymbol{x}$  belongs to class k;  $H_1$ :  $\boldsymbol{x}$  belongs to any class other than k }. Following the MAP principle, the MBR optimal decision should be based on the posterior  $p(H_b|\boldsymbol{x}): b = \{0,1\}$ , if the risk for  $H_0$  and  $H_1$  is symmetric. If the priors  $p(H_0)$  and  $p(H_1)$  are equal, we have:

$$p(H_b|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|H_b)}{p(\boldsymbol{x}|H_0) + p(\boldsymbol{x}|H_1)}.$$
(5)

Since  $p(H_0|\boldsymbol{x}) + p(H_1|\boldsymbol{x}) = 1$ , the decision can be simply made according to  $p(H_0|\boldsymbol{x})$ :

$$b^* = \begin{cases} 0 & if \ p(H_0 | \boldsymbol{x}) \ge 0.5 \\ 1 & if \ p(H_0 | \boldsymbol{x}) < 0.5. \end{cases}$$
(6)

In practice, by setting an appropriate threshold on  $p(H_0|\mathbf{x})$ , one can deal with different priors and risk on  $H_0$  and  $H_1$ . We highlight that for any class k, this threshold is only related to the prior and risk. This is important as it means that based on  $p(H_0|\mathbf{x})$ , MBR optimal decisions can be made *simultaneously* for all the classes by setting a global threshold. A simple case is to set the threshold to 0.5 when the risk is symmetric and the priors are equal. In summary,  $p(H_0|\mathbf{x})$  is MBR optimal for the SV task.

Note that when computing the posterior  $p(H_0|\mathbf{x})$ ,  $p(\mathbf{x}|H_0)$  is exactly the likelihood  $p_k(\mathbf{x})$ , and  $p(\mathbf{x}|H_1)$  summarizes the likelihood of all possible classes except the class k. In most cases, an SV system is required to deal with any unknown class, and so the class space is usually assumed to be continuous. To simplify the presentation, we will assume each class being uniquely represented by the mean vector  $\boldsymbol{\mu}$  and  $p(\boldsymbol{\mu})$  is continuous. In this case, the contribution of each class is infinitely small and so  $p(\boldsymbol{x}|H_1)$  is exactly the marginal distribution (or evidence)  $p(\boldsymbol{x}) = \int p(\boldsymbol{x}|\boldsymbol{\mu})p(\boldsymbol{\mu})d\boldsymbol{\mu}$ .<sup>[1]</sup> We therefore obtain the MBR optimal score for SV:

$$p(H_0|\boldsymbol{x}) = \frac{p_k(\boldsymbol{x})}{p_k(\boldsymbol{x}) + p(\boldsymbol{x})}.$$
(7)

# 2.1.3 Normalized likelihood

Note that for the SV task, according to Eq.(5), the posterior  $p(H_0|\boldsymbol{x})$  is determined by the ratio  $p(\boldsymbol{x}|H_0)/p(\boldsymbol{x}|H_1)$ , which is essentially the class-dependent likelihood  $p_k(\boldsymbol{x})$  normalized by the class-independent likelihood  $p(\boldsymbol{x})$ . We therefore define the normalized likelihood (NL) as:

$$NL(\boldsymbol{x}|k) = \frac{p(\boldsymbol{x}|H_0)}{p(\boldsymbol{x}|H_1)} = \frac{p_k(\boldsymbol{x})}{p(\boldsymbol{x})}.$$
(8)

Note that the NL is linked to the posterior  $p(H_0|\mathbf{x})$  by a monotone function:

<sup>&</sup>lt;sup>[1]</sup>One may argue that  $p(\mathbf{x})$  involves the quantity  $p_k(\mathbf{x})$ , and so is not accurately  $p(\mathbf{x}|H_1)$ . This is not true however, as the contribution of  $p_k(\mathbf{x})$  to  $p(\mathbf{x})$  is zero if  $p(\boldsymbol{\mu})$  is continuous. This also means that the likelihood that  $\mathbf{x}$  belongs to all classes equals to the likelihood that  $\mathbf{x}$  belongs to all classes other than k. Note that the prior  $p(\mu_k)$  is different from the prior  $p(H_0)$ :  $p(\mu_k)$  is the density that the class mean of a speaker is at  $\mu_k$ , while  $p(H_0)$  is the probability that a trial is positive, i.e., a genuine speaker.

$$NL(\boldsymbol{x}|k) = \frac{p(H_0|\boldsymbol{x})}{1 - p(H_0|\boldsymbol{x})}.$$
(9)

Since the posterior  $p(H_0|\mathbf{x})$  is MBR optimal for the SV task, the NL is also MBR optimal as a threshold on  $p(H_0|\mathbf{x})$  that leads to (global) MBR decisions can be simply transformed to a threshold on the NL, by which the same MBR decisions can be achieved. For example, the MBR decision is obtained when  $p(H_0|\mathbf{x}) = 0.5$  if the risk on  $H_0$  and  $H_1$  is equal, which is equal to say  $NL(\mathbf{x}|k) = 1.0$ , according to Eq.(9).

Interestingly, the NL score is also MBR optimal for the SI task. This is because the normalization term  $p(\mathbf{x})$  is the same for all classes in the SI task, so the decisions made based on the NL score is equal to those based on the likelihood  $p_k(\mathbf{x})$ . Since the likelihood is MBR optimal for the SI task, the NL score is MBR optimal for the SI task as well. We therefore conclude that the NL score is MBR optimal for both the SI and the SV tasks, under some appropriate assumptions. It should be noted that the NL form Eq. (8) is a high-level definition and it can be implemented in a flexible way. In particular,  $p_k(\mathbf{x})$  and  $p(\mathbf{x})$  can be any models that produce the class-dependent and class-independent likelihoods respectively.

Finally, NL is not new for speaker recognition. It is essentially the likelihood ratio (LR) that has been employed for many years since the GMM-UBM regime, where the score is computed by  $\frac{p_{GMM}(\boldsymbol{x})}{p_{UBM}(\boldsymbol{x})}$ . We use the term NL instead of LR in this paper in order to: (1) highlight the different roles of the numerator  $p_k(\boldsymbol{x})$  and the denominator  $p(\boldsymbol{x})$  in the ratio; (2) discriminate the normalization-style LR (used by NL) and the comparison-style LR, e.g., the one used by PLDA inference that compares the likelihoods that a group of samples are generated from the same and different classes.

#### 2.2 NL score with linear Gaussian model

Although the NL framework allows flexible models for the class-dependent and class-independent likelihoods, linear Gaussian model is the most attractive due to its simplicity. We derive the NL score with this model, for the case (1) the class means have been known and (2) the class means are unknown and have to be estimated from enrollment data.

## 2.2.1 Linear Gaussian model

We shall assume a simple linear Gaussian model for the speaker vectors that we will score:

$$p(\boldsymbol{\mu}) = N(\boldsymbol{\mu}; \mathbf{0}, \mathbf{I}\boldsymbol{\epsilon}^2) \tag{10}$$

$$p(\boldsymbol{x}|\boldsymbol{\mu}) = N(\boldsymbol{x};\boldsymbol{\mu},\sigma^2 \mathbf{I}),\tag{11}$$

where  $\boldsymbol{\mu} \in \mathbb{R}^D$  represents the means of classes and  $\boldsymbol{x} \in \mathbb{R}^D$  represents observations, and  $\boldsymbol{\epsilon}^2 \in (\mathbb{R}^+)^D$  and  $\sigma^2 \in \mathbb{R}^+$  represent the between-class and within-class variances respectively. Applied to speaker recognition,  $\boldsymbol{\epsilon}$  and  $\sigma$  represent the betweenspeaker and within-speaker variances respectively. We highlight that any linear Gaussian model can be transformed into this simple form (i.e., isotropic withinclass covariance and diagonal between-class covariance) by a linear transform such as full-dimensional linear discriminant analysis (LDA), and this linear transform will not change the identification and verification results as we will show in Section 2.3. Therefore, study with the simple form Eq. (10) and Eq. (11) is sufficient for us to understand the behavior of a general linear Gaussian model with complex covariance matrices.

With this model, it is easy to derive the marginal probability  $p(\mathbf{x})$  and the posterior probability  $p(\boldsymbol{\mu}|\mathbf{x})$  as follows [4]:

$$p(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{0}, \mathbf{I}(\boldsymbol{\epsilon}^2 + \sigma^2))$$
(12)

$$p(\boldsymbol{\mu}|\boldsymbol{x}) = N(\boldsymbol{\mu}; \frac{\boldsymbol{\epsilon}^2}{\boldsymbol{\epsilon}^2 + \sigma^2} \boldsymbol{x}, \mathbf{I} \frac{\sigma^2 \boldsymbol{\epsilon}^2}{\boldsymbol{\epsilon}^2 + \sigma^2}),$$
(13)

where all the operations between vectors are element-wised and appropriate dimension expansion has been assumed, e.g.,  $\boldsymbol{\epsilon}^2 + \sigma^2 = \boldsymbol{\epsilon}^2 + [\sigma^2, ..., \sigma^2]^T$ .

If the observations are more than one, the posterior probability has the form:

$$p(\boldsymbol{\mu}|\boldsymbol{x}_1,...,\boldsymbol{x}_n) = N(\boldsymbol{\mu}; \frac{n\boldsymbol{\epsilon}^2}{n\boldsymbol{\epsilon}^2 + \sigma^2} \bar{\boldsymbol{x}}, \mathbf{I} \frac{\sigma^2 \boldsymbol{\epsilon}^2}{n\boldsymbol{\epsilon}^2 + \sigma^2}),$$
(14)

where  $\bar{x}$  is the average of the observations. These equations will be extensively used in the following sections.

# 2.2.2 Case 1: class means are known

In this case, we assume that the class means are known. This is equivalent to say that each class is represented by infinite enrollment data.

# NL/Euclidean/Cosine scores for SI

For the SI tasks, decisions based on the NL score and the likelihood  $p_k(\boldsymbol{x})$  are the same and both are MBR optimal. With the linear Gaussian model, the likelihood is:

$$p_k(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}).$$
(15)

A simple rearrangement shows that:

$$\log p_k(\boldsymbol{x}) = -\frac{1}{2\sigma^2} ||\boldsymbol{x} - \boldsymbol{\mu}_k||^2 + const$$
(16)

Since the variance  $\sigma$  is the same for all classes, the MBR decision can be equally based on the Euclidean distance, e.g.,

$$s_e = ||\boldsymbol{x} - \boldsymbol{\mu}_k||^2, \tag{17}$$

where we use  $s_e$  to denote the score based on the Euclidean distance. In short, the Euclidean score is MBR optimal for the SI task when the class means are known.

Next, we will show that in a high-dimensional space, the Euclidean distance is well approximated by the cosine distance, under the linear Gaussian assumption.

First notice that the Gaussian annulus theorem [5] states that for a *d*-dimensional Gaussian distribution with the same variance  $\epsilon$  in each direction, nearly all the probability mass is concentrated in a thin annulus of width O(1) at radius  $\sqrt{d}\epsilon$ , as shown in Figure 1. This slightly anti-intuitive result indicates that in a high-dimensional space, most of the samples from a Gaussian tend to be in the same length. Rigid proof for this theorem can be found in [5]. Note that the distribution of real speaker vectors is not necessarily a perfect Gaussian; However, in most cases it can be well approximated by a Gaussian, especially when some normalization techniques are employed [23]. Therefore, the Gaussian annulus theorem can be readily used for speaker vectors.



Figure 1 Left: Gaussian annulus theorem [5]: for a *d*-dimensional multi-variant Gaussian with unit variance in all directions, for any  $\beta \leq \sqrt{d}$ , all but at most  $3e^{-c\beta^2}$  of the probability mass lies within the annulus  $\sqrt{d} - \beta \leq ||\mathbf{x}|| \leq \sqrt{d} + \beta$ , where *c* is a fixed positive constant. The color region shown in the figure represents the annulus. Rigid proof can be found in [5]. Right: The length distribution of samples from a 512-dimensional Gaussian distribution  $N(\mathbf{0}, \mathbf{I})$ .

Now we rewrite the Euclidean score as follows:

$$s_e = ||\boldsymbol{x}||^2 + ||\boldsymbol{\mu}_k||^2 - 2\cos(\boldsymbol{x}, \boldsymbol{\mu}_k)||\boldsymbol{x}|| \ ||\boldsymbol{\mu}_k||, \tag{18}$$

since  $||\boldsymbol{\mu}_k|| \approx \sqrt{d\epsilon}$ ,  $\cos(\boldsymbol{x}, \boldsymbol{\mu}_k)$  will be the only term that discriminates the probability that  $\boldsymbol{x}$  belongs to different class k. This leads to the cosine score:

$$s_c = \cos(\boldsymbol{x}, \boldsymbol{\mu}_k). \tag{19}$$

This result provides the theoretical support for the cosine score. It should be noted that this approximation is only valid for high-dimensional data, and the class means must be from a Gaussian with a zero mean. Therefore, data centralization is important for cosine scoring.

#### NL/Euclidean/Cosine scores for SV

For the SV task, the MBR optimal decision should be based on the NL score. With the linear Gaussian model, one can easily show that:

$$\log NL(\boldsymbol{x}|k) = -\frac{1}{2\sigma^2} ||\boldsymbol{x} - \boldsymbol{\mu}_k||^2 + \frac{1}{2} ||\frac{\boldsymbol{x}}{\sqrt{\boldsymbol{\epsilon}^2 + \sigma^2}}||^2 + const.$$
(20)

A simple rearrangement shows that:

$$\log NL(\boldsymbol{x}|k) = -\frac{1}{2\sigma^2} (||\boldsymbol{x}||^2 + ||\boldsymbol{\mu}_k||^2 - 2\cos(\boldsymbol{x}, \boldsymbol{\mu}_k)||\boldsymbol{x}|| ||\boldsymbol{\mu}_k||) + \frac{1}{2} ||\frac{\boldsymbol{x}}{\sqrt{\boldsymbol{\epsilon}^2 + \sigma^2}}||^2 + const$$

$$\propto -\left\{ ||\frac{\boldsymbol{\epsilon}}{\sqrt{\sigma^2 + \boldsymbol{\epsilon}^2}} \boldsymbol{x}||^2 + ||\boldsymbol{\mu}_k||^2 - 2\cos(\boldsymbol{x}, \boldsymbol{\mu}_k)||\boldsymbol{x}|| ||\boldsymbol{\mu}_k|| \right\}.$$
(21)

It can be seen that if the within-class variance  $\sigma^2$  is significantly larger than the between-class variance  $\epsilon^2$  (we refer to element-based comparison here and after), the log NL will significantly depart from the Euclidean distance, but more closely related to the cosine distance. Essentially, if we admit that both  $||\boldsymbol{x}||^2$  and  $||\boldsymbol{\mu}_k||^2$  tend to be constant due to the Gaussian annulus theorem, the cosine score will be a good approximation for the optimal log NL. Conversely, if the between-class variance  $\epsilon^2$  is sufficient larger than the within-class variance  $\sigma^2$ , it can be well approximated by the Euclidean score.

## 2.2.3 Case 2: class means are unknown

In the pervious section, we have supposed that the class means are known precisely. In real scenarios, however, this is not possible. We usually have only a few enrollment samples (e.g., less than 3) to represent the class, and the SI or SV evaluation should be based on these representative samples. In this case, the class means are unknown and have to be estimated from the enrollment data, leading to uncertainty that must be taken into account during scoring.

#### NL/Euclidean/Cosine scores for SI

Firstly consider the MBR optimal decision for SI. As in the known-mean scenario, we compute the likelihood for the k-th class:

$$p_k(\boldsymbol{x}) = N(\boldsymbol{x}; \boldsymbol{\mu}_k, \sigma^2 \mathbf{I}).$$
<sup>(22)</sup>

An important difference here is that  $\boldsymbol{\mu}_k$  is unknown, and so has to be estimated from the enrollment samples belonging to the same class. Denoting these samples by  $\boldsymbol{x}_1^k, \dots \boldsymbol{x}_{n_k}^k$  and their average by  $\bar{\boldsymbol{x}}_k$ , we have the posterior probability for the class mean  $\boldsymbol{\mu}_k$ , according to Eq. (14):

$$p(\boldsymbol{\mu}_k | \boldsymbol{x}_1^k, \dots \boldsymbol{x}_{n_k}^k) = N(\boldsymbol{\mu}_k; \frac{n_k \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2} \bar{\boldsymbol{x}}_k, \mathbf{I} \frac{\sigma \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2}).$$
(23)

The likelihood  $p_k(\boldsymbol{x})$  can therefore be computed by marginalizing over  $\boldsymbol{\mu}_k$ , according to this posterior. Following Eq.(12), we have:

$$p_{k}(\boldsymbol{x}) = p(\boldsymbol{x}|\boldsymbol{x}_{1}^{k},...,\boldsymbol{x}_{n_{k}}^{k})$$

$$= \int p(\boldsymbol{x}|\boldsymbol{\mu}_{k})p(\boldsymbol{\mu}_{k}|\boldsymbol{x}_{1}^{k},...\boldsymbol{x}_{n_{k}}^{k})d\boldsymbol{\mu}_{k}$$

$$= N(\boldsymbol{x};\frac{n_{k}\boldsymbol{\epsilon}^{2}}{n_{k}\boldsymbol{\epsilon}^{2}+\sigma^{2}}\bar{\boldsymbol{x}}_{k},(\sigma^{2}+\mathbf{I}\frac{\sigma\boldsymbol{\epsilon}^{2}}{n_{k}\boldsymbol{\epsilon}^{2}+\sigma^{2}})).$$
(24)

Note that with the class mean uncertainty, the Euclidean score is not MBR optimal anymore. If the number of enrollment observations are the same for all classes, the likelihood is exclusively determined by the class mean  $\boldsymbol{\mu}_k$ . In this case, an amended version of the Euclidean score is optimal, where the class mean is computed by  $\frac{n_k \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2} \boldsymbol{\mu}_k$ . Note that the scale  $\frac{n_k \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2}$  has been applied to compensate for the uncertainty of the maximum-likelihood mean estimation  $\boldsymbol{\mu}_k$ . Intuitively, a smaller n or a larger  $\sigma^2/\boldsymbol{\epsilon}^2$  lead to more uncertainty, so the compensation term will be more significant. With more enrollment samples, the compensation term will converge to one, and the standard Euclidean score is recovered.

Another observation is that the scale compensation on  $\mu_k$  does not change its direction. This implies that the cosine score does not need any amendment to account for the uncertainty. However, it does not mean that the cosine score is not impacted by the class mean uncertainty; it just means that the cosine score is not impacted as much as the Euclidean score.

# NL/Euclidean/Cosine scores for SV

Now we normalize the score  $p_k(\boldsymbol{x})$  to make it suitable for SV, by introducing a normalization term  $p(\boldsymbol{x})$ :

$$NL(\boldsymbol{x}|k) = \frac{p(\boldsymbol{x}|\boldsymbol{x}_1, \dots, \boldsymbol{x}_{n_k})}{p(\boldsymbol{x})}.$$
(25)

Note that the normalization term  $p(\mathbf{x})$  is not impacted by the mean uncertainty, and therefore remains the same value as in the known-mean scenario. A simple computation shows that:

$$\log NL(\boldsymbol{x}|k) \propto -||\frac{\boldsymbol{x} - \tilde{\boldsymbol{\mu}}_k}{\sqrt{\sigma^2 + \frac{\boldsymbol{\epsilon}^2 \sigma^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2}}}||^2 + ||\frac{\boldsymbol{x}}{\sqrt{\boldsymbol{\epsilon}^2 + \sigma^2}}||^2,$$
(26)

where we have defined:

$$\tilde{\boldsymbol{\mu}}_{k} = \frac{n_{k} \boldsymbol{\epsilon}^{2}}{n_{k} \boldsymbol{\epsilon}^{2} + \sigma^{2}} \bar{\boldsymbol{x}}_{k}.$$
(27)

To compare with the Euclidean score and the cosine score, Eq. (25) can be reformulated to:

$$\log NL(\boldsymbol{x}|k) \propto -\left\{\frac{n_k \boldsymbol{\epsilon}^4}{(\sigma^2 + \boldsymbol{\epsilon}^2)(n_k \boldsymbol{\epsilon}^2 + \sigma^2)} ||\boldsymbol{x}||^2 + ||\boldsymbol{\tilde{\mu}}_k||^2 - 2\cos(\boldsymbol{x}, \boldsymbol{\tilde{\mu}}_k)||\boldsymbol{x}|| ||\boldsymbol{\mu}_k||\right\}.$$
(28)

It can be seen that if the between-class variance  $\boldsymbol{\epsilon}$  is significantly smaller than the within-class variance  $\sigma$ , the first two terms on the right hand side of Eq. (28) tend to be small and log NL can be approximated by the cosine score. On the opposite, if the between-class variance  $\boldsymbol{\epsilon}$  is significantly larger than the within-class variance  $\sigma$ , the amended Euclidean score will be a good approximation. Finally, if  $n_k$  is sufficiently large, Eq. (28) will fall back to Eq. (21) of the know-mean case.

# 2.3 Remarks on properties of NL score Remark 1: Equivalent to PLDA LR

The NL score based on the linear Gaussian model and unknown class means is equivalent to the PLDA LR [28, 46]. PLDA assumes the same linear Gaussian model, but uses the following likelihood ratio as the score:

$$LR_{PLDA}(\boldsymbol{x}) = \frac{p(\boldsymbol{x}, \boldsymbol{x}_1, ..., \boldsymbol{x}_n \text{ from the same class})}{p(\boldsymbol{x} \text{ from a unique class})p(\boldsymbol{x}_1, ..., \boldsymbol{x}_n \text{ from a unique class})}$$

Note that this likelihood ratio is different from the likelihood ratio of the NL score in Eq. (8). The PLDA LR can be formally represented by:

$$LR_{PLDA}(\boldsymbol{x}) = \frac{p(\boldsymbol{x}, \boldsymbol{x}_1, \dots, \boldsymbol{x}_n)}{p(\boldsymbol{x})p(\boldsymbol{x}_1, \dots, \boldsymbol{x}_n)}$$
(29)

where  $p(\boldsymbol{x}_1, ..., \boldsymbol{x}_n)$  denotes the probability that  $\boldsymbol{x}_1, ..., \boldsymbol{x}_n$  belong to the same but an unknown class. In principle, this quantity can be computed by marginalizing over the class mean:

$$p(\boldsymbol{x}_1,...,\boldsymbol{x}_n) = \int p(\boldsymbol{x}_1,..,\boldsymbol{x}_n | \boldsymbol{\mu}) p(\boldsymbol{\mu}) \mathrm{d}\boldsymbol{\mu}.$$
(30)

A simple re-arrangement shows that:

$$LR_{PLDA}(\boldsymbol{x}) = \frac{p(\boldsymbol{x}|\boldsymbol{x}_1, ..., \boldsymbol{x}_n)}{p(\boldsymbol{x})} = \frac{\int p(\boldsymbol{x}|\boldsymbol{\mu})p(\boldsymbol{\mu}|\boldsymbol{x}_1, ..., \boldsymbol{x}_n) \mathrm{d}\boldsymbol{\mu}}{p(\boldsymbol{x})},$$
(31)

where we have divided the numerator  $p(\boldsymbol{x}, \boldsymbol{x}_1, ..., \boldsymbol{x}_n)$  by  $p(\boldsymbol{x}_1, ..., \boldsymbol{x}_n)$ , which converts the marginal distribution  $p(\boldsymbol{x}, \boldsymbol{x}_1, ..., \boldsymbol{x}_n)$  to the conditional distribution  $p(\boldsymbol{x}|\boldsymbol{x}_1, ..., \boldsymbol{x}_n)$ . By this change, the numerator is the likelihood of  $\boldsymbol{x}$  belonging to the class represented by  $\boldsymbol{x}_1, ..., \boldsymbol{x}_n$ , and the denominator is the likelihood  $\boldsymbol{x}$  belonging to any class. This is exactly the normalized likelihood in Eq.(25). We therefore

conclude that the PLDA LR is an NL where the underlying probabilistic model is linear Gaussian and the class means are estimated from finite enrollment data. Since the NL score is MBR optimal for both SI and SV tasks, an immediate conclusion is that the PLDA LR is also MBR optimal for the two tasks. Note that the NL form of the PLDA LR was discussed by McCree et. al. [6, 42].

Compared to PLDA LR, NL possesses some attractive properties and brings some interesting merits. A particular merit is that NL decouples the score computation into three steps: posterior computation based on enrollment data, likelihood computation for the test data based on the posterior, and normalization based on a global model. This offers an interesting correspondence between the scoring model and the scoring process. We therefore can investigate the behavior of each component and design fine-grained treatment for real-life imperfection, e.g., the enrollment-test mismatch that will be presented in Section 3.7.

## **Remark 2: Invariance with invertible transform**

Suppose an invertible transform g on  $\boldsymbol{x}$ , and the probabilities on  $\boldsymbol{x}$  and  $g(\boldsymbol{x})$  are p and p' respectively. According to the principle of distribution transformation for continuous variables [50], p and p' has the following relation:

$$p'(g(\boldsymbol{x})) = p(\boldsymbol{x}) |\det \frac{\partial g^{-1}(\boldsymbol{x})}{\partial \boldsymbol{x}}|, \qquad (32)$$

where the second term is the absolute value of the determinant of the Jacobian matrix of  $g^{-1}$ , the inverse transform of g. This term reflects the change of the volume with the transform, and is often called the entropy term and denoted by  $J(\boldsymbol{x})$ .

For the marginal distribution  $p(\boldsymbol{x}_1,...,\boldsymbol{x}_n)$  where  $\boldsymbol{x}_1,...,\boldsymbol{x}_n$  are drawn from the same but an unknown class, one can compute the distribution by:

$$p(\boldsymbol{x}_1, ..., \boldsymbol{x}_n) = \int \prod_{i=1}^n p(\boldsymbol{x}_i | \boldsymbol{\mu}) p(\boldsymbol{\mu}) d\boldsymbol{\mu}$$
$$\approx \sum_j \prod_{i=1}^n p(\boldsymbol{x}_i | \boldsymbol{\mu}_j) p(\boldsymbol{\mu}_j) \Delta(\boldsymbol{\mu}_j)$$
(33)

where we have divide the  $\boldsymbol{\mu}$  space into a large amount of small areas  $\{\Delta(\boldsymbol{\mu}_j)\}$  with centers  $\{\boldsymbol{\mu}_j\}$ . The approximation will approach to be accurate when the number of small areas is infinite. With the transform g, we have:

$$p'(g(\boldsymbol{x}_1), ..., g(\boldsymbol{x}_n)) \approx \sum_j \prod_{i=1}^n p'(g(\boldsymbol{x}_i) | \boldsymbol{\mu}_j^g) p'(g(\boldsymbol{\mu}_j)) \Delta(g(\boldsymbol{\mu}_j))$$
(34)

where  $\boldsymbol{\mu}_{j}^{g}$  represents the mean of the class centered at  $\boldsymbol{\mu}_{j}$  after the transform. Moreover, the transform g does not change the probability within  $\Delta(\boldsymbol{\mu}_{j})$ , which means:

$$p(\boldsymbol{\mu}_j)\Delta(\boldsymbol{\mu}_j) = p'(g(\boldsymbol{\mu}_j))\Delta(g(\boldsymbol{\mu}_j)).$$
(35)

Putting all the pieces together, we have:

$$p'(g(\boldsymbol{x}_1),...,g(\boldsymbol{x}_n)) \approx \sum_j \prod_{i=1}^n p'(g(\boldsymbol{x}_i)|\boldsymbol{\mu}_j^g))p(\boldsymbol{\mu}_j)\Delta(\boldsymbol{\mu}_i)$$
$$= \sum_j \prod_{i=1}^n J(\boldsymbol{x}_i) \prod_{i=1}^n p(\boldsymbol{x}_i|\boldsymbol{\mu}_j)p(\boldsymbol{\mu}_j)\Delta(\boldsymbol{\mu}_i),$$
(36)

where we have applied the rule of the distribution transform shown in Eq. (32). Let the size of  $\{\Delta(\boldsymbol{\mu}_j)\}$  to be infinite, we have the marginal distribution in the space induced by transform g:

$$p'(g(\boldsymbol{x}_1), ..., g(\boldsymbol{x}_n)) = \prod_{i=1}^n J(\boldsymbol{x}_i) p(\boldsymbol{x}_1, ..., \boldsymbol{x}_n).$$
(37)

Substituting back to the NL score, we obtain the invariance of the NL score under an invertible transform:

$$NL(g(\boldsymbol{x})|g(\boldsymbol{x}_{1}),...,g(\boldsymbol{x}_{n_{k}})) = \frac{p'(g(\boldsymbol{x}),g(\boldsymbol{x}_{1}),...,g(\boldsymbol{x}_{n_{k}}))}{p'(g(\boldsymbol{x}))p'(g(\boldsymbol{x}_{1}),...,g(\boldsymbol{x}_{n_{k}}))}$$
$$= \frac{J(\boldsymbol{x})\prod_{i=1}^{n}J(\boldsymbol{x}_{i})p(\boldsymbol{x},\boldsymbol{x}_{1},...,\boldsymbol{x}_{n_{k}})}{\{J(\boldsymbol{x})p(\boldsymbol{x})\}\{\prod_{i=1}^{n}J(\boldsymbol{x}_{i})p(\boldsymbol{x}_{1},...,\boldsymbol{x}_{n_{k}})\}}$$
$$= NL(\boldsymbol{x}|\boldsymbol{x}_{1},...,\boldsymbol{x}_{n_{k}})$$
(38)

where we have employed the PLDA LR form to represent the NL score.

The above derivation indicates that the NL score can be computed in a transformed space induced by an invertible transform. Among all the possible invertible transforms, the full-dimension LDA is particularly attractive. It can simultaneously diagonalize the within-class and between-class covariances and regulate the withinclass covariance to be identity. We therefore do not need consider the general form of distributions when investigating the properties of the NL score, instead just focusing on the simple form with diagonal covariances, as we did in the previous sections.

## **Remark 3: Dimensionality is important**

Let's investigate the benefit of a high-dimensional space. It has been shown [5] that the distance of two random samples from a *n*-dimensional Gaussian with variance  $\epsilon^2$  in all directions has a large probability to be:

$$||\boldsymbol{x} - \boldsymbol{y}|| = \sqrt{2d\epsilon} \pm O(1). \tag{39}$$

Consider the class means are random samples of a Gaussian with variance  $\epsilon^2$ , and each class is a Gaussian with variance  $\sigma^2$ . Due to the Gaussian annulus theorem, the samples of each class will concentrate in the annulus of radius  $\sqrt{d\sigma}$ . Since the

distance of two class means has a large probability to be  $\sqrt{2d\epsilon}$ , it is easy to conclude that if  $2\sigma < \epsilon$ , there will be a large probability that most of the classes are well separated.

More careful analysis shows a better bound. Considering two samples from two different classes respectively, it shows that their distance tend to be  $\sqrt{\Delta^2 + 2\sigma^2 d \pm O(\sqrt{d}\sigma)}$ , where  $\Delta$  is the distance of these two classes, and  $\sigma^2$  is the variance of each class [5]. Since the samples from the same class tends to be  $\sqrt{2d\sigma}$ , one can show if  $\Delta^2 \ge O(\sqrt{d\sigma})$ , there will be a large probability to identify if two samples are from the same class or different classes. If the class means are sampled from a Gaussian with variance  $\epsilon$ , we will have  $\Delta^2 \approx 2\epsilon^2 d$ . One can easily derive that if  $\sigma^2 \le O(\epsilon^4 d)$ , sample pairs from two classes can be well differentiated from sample pairs from the same class. Note the condition depends on d, which means that with a higher dimension, classes with larger variances can be separated with a large probability. In other words, classes in higher dimensional space tend to be more separable.

# **Remark 4: Direction is important**

Another interesting property of a high dimension space is that most of the volume of a unit ball is concentrated near its "equator" [5], as shown in Figure 2. More precisely, for any unit-length vector v defining the "north", most of the volume of the unit ball lies in the thin slab of points whose dot-product with v has magnitude  $O(\frac{1}{\sqrt{d}})$  [5].



An immediate conclusion is that for any sample from a Gaussian, it is orthogonal to most of other samples from the same Gaussian. This is evident if we note that the dot product of any two samples tend to be  $1/\sqrt{d}$ , which approaches to zero with a large d. Combining the Gaussian annulus theorem, we can see that samples of a high-dimensional Gaussian are mostly scattered across direction rather than length. In other words, **direction is more important than magnitude in a high dimensional space.** In fact, the importance of direction in high-dimensional space has been noticed by researchers in various domains. For example, it is well-known that the cosine distance is a better metric compared to the Euclidean distance in text analysis and information retrieval [13, 52, 67]. The same observation was also reported in speaker recognition [16, 31].

It is worth noting that all the above conclusions are based on Gaussian distributions. If the data itself is spherical in nature, a directional distribution will be naturally preferred, for example the Von Mises-Fisher (VMF) distribution. More information about directional distributions can be found in [41, 56].

# **3** Results

In this section, we will discuss the application of the NL score in practical speaker recognition systems. For simplicity, we only focus on the NL score based on the linear Gaussian model. The main purpose is to investigate the behavior of the NL score. Theoretically, NL scoring is MBR optimal if the data satisfy the model assumption, and real-life imperfection is essentially represented by the mismatch between the distributions that the model assumes and the data exhibit. We therefore conduct the investigation by simulating this mismatch, one type per experiment. Note that all the EER/IDR results reported in this section are based on the NL score.

In order to reflect the behavior of the NL score in real-life systems, we need consider: (1) The true configuration of practical speaker vectors, including the number of dimensions and classes, the range of the between-class and within-class variances. These configurations will provide information about the operation point of the NL scoring, by which we can obtain the expected performance of a speaker recognition system if the linear Gaussian assumption is satisfied. (2) The deviation of the distribution of practical speaker vectors from the linear Gaussian assumption, in particular the potential problem of non-homogeneity and non-Gaussianlity. The former concerns how different speakers differ from each other, and the latter concerns how the between-class distribution of the speaker means and the within-class distributions of individual speakers deviate from Gaussian. By these information, we can estimate how much performance loss would be expected in practical systems with the NL scoring.

# 3.1 Baseline systems

**Data** We use the VoxCeleb [14, 43] dataset to build an x-vector system and an i-vector system. The entire database consists of *VoxCeleb1* and *VoxCeleb2*. All the speech signals were collected from open-source media channels and therefore involve rich variations in channel, style, and ambient noise. The entire dataset contains 2000+ hours of speech signals from 7000+ speakers. Data augmentation was applied to improve robustness, with the MUSAN corpus [54] used to generate noisy utterances, and the room impulse responses (RIRS) corpus [33] was used to generate reverberant utterances.

**x-vector system**: The x-vector frontend was created using the Kaldi toolkit [45], following the SITW recipe. The acoustic features are 40-dimensional Fbanks. The main architecture contains three components. The first component is the feature-learning component, which involves 5 time-delay (TD) layers to learn frame-level speaker features. The slicing parameters for these 5 TD layers are:  $\{t-2, t-1, t, t+1, t+2\}$ ,  $\{t-2, t, t+2\}$ ,  $\{t-3, t, t+3\}$ ,  $\{t\}$ ,  $\{t\}$ . The second component is the statistical pooling component, which computes the mean and standard deviation of the frame-level features from a speech segment. The final one is the speaker-classification component, which discriminates between different speakers. This component has 2 full-connection (FC) layers and the size of its output is 7,185, corresponding to the number of speakers in the training set. Once trained, the 512-dimensional activations of the penultimate FC layer are read out as an x-vector.

**i-vector system**: The i-vector frontend was built with the Kaldi toolkit [45], following the SITW recipe as well. The raw features involve 24-dimensional MFCCs plus the log energy, augmented by first- and second-order derivatives, resulting in a 75-dimensional feature vector. This feature is used by the i-vector model. The universal background model (UBM) consists of 2,048 Gaussian components, and the dimensionality of the the i-vectors is set to be 400.

# 3.2 Statistics of x-vectors and i-vectors

We first look at the properties of different types of speaker vectors. To ensure sufficient statistical strength, we choose 4000 speakers with sufficient utterances from the VoxCeleb training data. The number of utterances per speaker in this set is 45 in average, and the minimum and maximum values are 10 and 438, respectively. All the data are preprocessed by a full-dimension LDA, by which the *accumulated* within-class covariance is normalized to be an identify matrix, and the between-class covariance becomes diagonal. Note that the full-dimension LDA does not change the NL scores, but the simplified covariance structure makes the analysis easier. We compute a number of statistics, regarding the homogeneity (i.e., if all classes share the same covariance) and Gaussianality of the within-class and between-class distributions.

- PC direction STD for homogeneity. This tests if the covariance matrices of all the speakers have the same direction. After PCA, the first principle component (PC1) of all the speakers are selected and its mean over the speakers is computed. The cosine distance between the PC1s of individual speakers and the mean PC1 is computed. The STD of these cosine scores is used as the measure to test the PC1 direction variance. The same computation is conducted on all PCs. In this experiment, we report the direction variance on PC1 and PC2, and the averaged direction variance on the first 10 PCs.
- PC shape STD for homogeneity. Using PC1 as an example, the coefficients (eigenvalues) of the covariance matrices of all the speakers on the first PC are calculated, and the STD of these coefficients over all speakers is computed. The same computation is performed on all the PCs. Since the coefficient on each PC determines the spreading of the samplings on this direction, the coefficients on all the PCs determine the shape of the speaker distribution. The STD of these coefficients over all speakers then test if the distributions of all speakers have the same shape (regardless of the direction), hence being noted as PC shape STD. We report the PC shape STD on PC1 and PC2, and the averaged PC shape STD on the first 10 PCs.
- Averaged PC kurtosis for Gaussianality. On each PC direction, we compute the kurtosis for each speaker, and then compute the mean of the kurtosis over all the speakers. The averaged kurtosis over the first 10 PCs is reported.
- Averaged PC skewness for Gaussianality. On each PC direction, we compute the skewness for each speaker, and then compute the mean of the kurtosis over all the speakers. The averaged skewness over the first 10 PCs is reported.
- Between-class kurtosis and skewness. The kurtosis and skewness of the class means, computed on each dimension, and then are averaged.

To have a comparison with the ideal case where the data are truly linear Gaussian, synthesis datasets are constructed for the x-vectors and i-vectors respectively. We first sample the same number of classes (4000) using the same between-class covariance of the true speaker vectors. For each class, we sample the same number of samples of that class in the real data, using the same within-class covariance (which is 1.0 in the LDA space). From these synthesis data, we compute the same statistics as the real speaker vectors. These values can be used to evaluate how the real dataset departs from a perfect linear Gaussian dataset.

The results are shown in Table 1. It can be seen that the real speaker vectors exhibit clear non-homogeneity and non-Gaussianality. For non-homogeneity, it looks like the most variance lies on the shape rather than the direction of the withinclass distributions. Moreover, the x-vectors and i-vectors show similar shape and direction variances, which means that these two types of speaker vectors are not much different in terms of non-homogeneity.

For non-Gaussianality, both the x-vectors and i-vectors are clearly non-Gaussian, in terms of both between-class and within-class distributions. Specifically, it seems that the most significant difference between x-vectors and i-vectors is that the kurtosis of the within-class distribution is much higher with the x-vectors, and the large positive value suggests that the x-vectors mostly concentrate on the class means.

As for the between-class distribution, it seems that for both the x-vectors and i-vectors, the distribution is Gaussian, and the difference between the two kinds of speaker vectors is not substantial.

We also compute the EER and IDR results with the NL scoring. In this test, one sample from each class is used for enrollment and one sample is used for test. To ensure statistical significance, we run the test 500 times, and report the EER and IDR results as well as the variation. The results are shown in the bottom of Table 1. It can be seen that if the data are truly linear Gaussian, the NL scoring ensures a very high performance. This performance is an upper bound that the NL scoring can achieve. In real-life situations, this upper bound is hard to reach, due to the non-homogeneity and non-Gaussianality of the data in nature, as well as the complexity associated with domain and condition mismatch.

In the reset of this section, we will conduct a series of simulation experiments, to study the impact of various factors related to the real-life imperfection. We hope this analysis will help identify the key factors that should be cared when designing a practical speaker recognition system. Due to the superior performance of x-vectors, our simulation will be based on the x-vector configuration. The NL score is used in all the following experiments.

#### 3.3 Problem associated with non-Gaussianality

In the previous section, we have found that x-vectors are highly non-Gaussian, particularly in terms of kurtosis of the within-class distributions. We perform a simulation experiment to investigate the impact of a high kurtosis. We use the Laplace distribution whose excessive kurtosis is 3. This is not as high as the x-vectors showed, but at least higher than the value of a Gaussian. The experiment is based on the configuration of the x-vectors derived from VoxCeleb (in the LDA space). We sample 600 classes following the same between-class distribution as the x-vectors.

		x-vector		i-vector	
		True Data	Synthesis Data	True Data	Synthesis Data
Within Class	PC1 dir. STD	0.095	0.045	0.060	0.045
	PC2 dir. STD	0.077	0.045	0.063	0.055
	Avg PC dir. STD	0.084	0.045	0.055	0.055
	PC1 shape STD	11.2	1.80	11.5	1.94
	PC2 shape STD	6.83	1.82	7.95	1.95
	Avg PC shape STD	5.53	1.85	5.48	1.98
	PC Kurtosis	19.15	0.150	1.940	0.152
	PC Skewness	0.837	0.040	0.370	0.041
Between Class	PC Kurtosis	0.390	0.010	0.358	0.002
	PC Skewness	0.065	0.001	0.043	0.001
	EER% (STD)	2.77 (0.1)	0.0 (0.0)	3.22 (0.2)	0.003 (0.008)
	IDR% (STD)	85.51(0.7)	100.0 (0.0)	75.08 (0.5)	100.0(0.0)

Table 1 Statistics of x-vectors and i-vectors, and the corresponding synthesis data.

For each class, we sample one sample for enrollment and three samples for test, from either a Gaussian or a Laplace distribution. In the NL scoring, we assume all the data are generated from Gaussian, and use the within-class variance that is used to generate the data. Each test repeats 500 times, and the averaged EER and IDR are reported, plus their variations. The results with different within-class variances are shown in Figure 3. It can be seen that the incorrect Laplace distribution indeed detriments the performance, especially in terms of EER. With a high within-class variance, the Laplace distribution seems hurt the IDR performance not much, which may be attributed to the fact that the Laplace distribution is more concentrated than the Gaussian. We conjecture that a larger kurtosis will lead to more severe performance reduction.



Researchers have noticed the problem associated with non-Gaussianality. Various non-linear transforms have been proposed, for example the radial Gaussianization (RG) [40] and the simple length normalization [23]. For x-vectors, involving Gaussian constraints in the training objective of the x-vector extractor may improve the Gaussianality [8, 36, 37]. Variational auto-encoders (VAE) and normalization flows [58, 65, 69] were also employed to improve Gaussianlity of x-vectors. An-

other line of research employs a non-Gaussian model, with the hope to handle non-Gaussian data in practical situations [31] .

#### 3.4 Problem associated with non-homogeneity

The non-homogeneity is caused by the variation of individual classes. The results in Table 1 show that this variation is largely related to the shape rather than the direction of the distributions of individual classes. We therefore focus on the impact of the variation of within-class variances, i.e., variance's variation.

We perform a simulation test, by imitating the between-class and (accumulated) within-class variance (in the LDA space) of x-vectors derived from VoxCeleb. A noise will be added to the variance of each individual class, to simulate the non-homogeneity.

Specifically, we sample 600 classes according to the between-class distribution of the x-vectors. For each class, we sample one sample for enrollment, and three samples for test. The variance of each class will be modified during the sampling by adding a noise  $\xi$ , but the same within-class distribution is used when sampling the enrollment and test data for that class. More specially, when sampling data for a particular class, a random noise  $\xi$  is added to the STD of the within-class distribution (1.0 in our test). Note that when the within-class variance is smaller than 0.1 after adding the noise, we will keep the variance to be 0.1. In our experiment, we test the impact of different levels of non-homogeneity, by varying the STD value of the added noise from 0.1 to 3.0. Therefore, the final within-class variance is max(0.1, 1.0 +  $\xi$ ), where  $\xi \sim N(0, \omega)$  and  $\omega$  changes from 0.1 to 3.0. For each  $\omega$ , the test runs 500 rounds and the mean and variation of the EER and IDR results are reported, on the SV and SI tasks respectively.

Since adding noise to the within-class variance of individual classes will change the accumulated within-class variance, the original configuration (within-class variance = 1.0) is not correct for NL scoring. We generate 200 samples for each class with exactly the same variance (after adding noise) of each class, and then compute the accumulated within-class variance using these samples. This accumulated within-class variance is used for computing the NL score of the non-homogeneous dataset. Additionally, we also generate a homogeneous dataset, where all the classes are generated using the same accumulated within-class variance used within-class variance is used for computing the NL score of the non-homogeneous dataset. Additionally, we also generate a homogeneous dataset, where all the classes are generated using the same accumulated within-class variance used when sampling the non-homogeneous dataset. This will be used as the homogeneous reference for the comparative analysis.

The results are shown in Figure 4. It can be seen that the non-homogeneous data generally achieve worse performance compared to the homogeneous data, in terms of both EER and IDR. An exception is that when the noise STD is 3.0, the IDR performance of the non-homogeneous data is better than the homogeneous data. This is attributed to the fact that according to our sampling scheme, a large portion of the within-class variances collapse to 0.1 when the non-homogeneous level is high, leading to a subset of classes whose within-class distributions are not only homogeneous but also compact. This is not a real-life situation. Note that the compact data suffer from biased within-class distribution, hence a worse EER.

The research on non-homogenity is far from extensive. The central loss that imposes the same Gaussian constraint for individual classes may improve homogeneity [8, 36, 37]. Recently we presented a deep normalization approach [9] based on



normalization flows [19, 20, 32]. This approach intends to regulate individual classes into a standard Gussian by a deep neural net, and has achieved promising results with x-vectors.

# 3.5 Problem associated with training-deployment domain mismatch

Besides the break of the linear Gaussian assumption, NL also suffers from incorrect configurations, i.e., using incorrect between-class and/or within-class covariances when computing the NL score. In practice, this often happens when the NL parameters are estimated in one domain (training phase), but are used in another domain (deployment phase). We will investigate the factors that mostly impact the NL scoring under the training-deployment mismatch by simulation experiments.

#### 3.5.1 Statistical analysis

To understand what has been changed from one domain to another, we compare the distributional properties of the x-vectors derived from VoxCeleb and another dataset, CNCeleb [21]. The two datasets are in different languages and with different genres, so can represent two domains.

Firstly, we compute the between-class and within-class covariances of the two datasets, shown in Figure 5 and Figure 6 respectively. Then we use VoxCeleb to train an LDA, and apply it to transform data of both VoxCeleb and CNCeleb. Note that LDA does not change the NL behavior, but can regularize the data to a simple distribution, making the comparison of the two datasets easier.

The between-class and within-class covariances of VoxCeleb and CNCeleb after the LDA transform are shown in Figure 7 and Figure 8 respectively. It can be seen that the LDA trained on VoxCeleb can largely diagonalize the between-class and within-class covariances of CNCeleb. This is a nice property and suggests that the directions of the distribution of the class means (related to between-class covariance) and the accumulated distribution of individual classes (related to within-class covariance) do not change significantly from VoxCeleb to CNCeleb. Note that in Table 1, we have shown that the directional variance of individual classes is small. Therefore, we conclude that the directions of both the between-class distribution and the individual within-class distributions do not change much from VoxCeleb to CNCeleb.

However, the diagonal elements of the between-class covariance do change significantly. Specifically, for VoxCeleb, most of the variance is distributed over the first several dimensions; for CNCeleb, however, the distribution tends to be uniform. For the within-class covariance, the diagonal elements remain equally distributed, but the value of each element has changed significantly from VoxCeleb to CNCeleb (1.0 for VoxCeleb vs 3.9 for CNCeleb).





More quantitative analysis are shown in Table 2, where we have shown the statistics of the x-vectors transformed by two LDAs, trained on VoxCeleb and CNCeleb respectively. For the within-class (WC) and between-class (BC) covariances in the LDA space, we compute the mean/variance of the diagonal elements as well as the proportion of the values of the diagonal elements (concentration factor). A key observation is that the concentration factors of the between-class and within-class







covariances are relatively large, by applying the LDA learned from either Vox-Celeb and CNCeleb. It double confirms that the directions of the within-class and between-class distributions do not change much from one domain to another.

Moreover, the variation of the diagonal elements of the within-class covariance is relatively small for both datasets, though the mean of the diagonal elements is different for different datasets. This indicates that the within-class distribution has changed from one dataset to another, but the change can be simply compensated by a global scale factor on all the dimensions. The between-class covariance shows different properties. The variation of the diagonal elements is much higher on the data where the LDA is trained, indicating that the between-class distribution has been changed significantly from one dataset to another.

 Table 2
 Statistics of x-vectors in VoxCeleb and CNCeleb, after applying LDAs trained on each of them.

		VoxCeleb	CNCeleb
LDA by VoxCeleb	WC Mean	1.000	4.094
	WC Var	0.000	0.138
	WC Concentration	1.000	0.971
	BC Mean	0.764	1.142
	BC Var	11.22	2.870
	BC Concentration	1.000	0.971
LDA by CNCeleb	WC Mean	0.534	1.000
	WC Var	0.019	0.000
	WC Concentration	0.930	1.000
	BC Mean	0.721	0.507
	BC Var	0.722	1.103
	BC Concentration	0.930	1.000

# 3.5.2 Simulation results

We perform a simulation experiment to test the degradation that the domain mismatch may cause. Again, the simulation is based on the configuration of the xvectors derived from VoxCeleb, but the between-class and within-class variances will be changed to smaller or larger when sampling the enrollment/test data. For a better comparison, the within-class variance of the baseline (without domain mismatch) is set to be 2.0. In each test, we sample 600 classes, and for each class, we sample one sample for enrollment and three samples for test. We run each test for 500 rounds, and report the averaged EER and IDR, plus the variations on them.

The first experiment simulates the impact with incorrect between-class variances. For that purpose, we define a distortion factor  $\alpha$ , and multiply the original betweenclass variances (on all dimension) by  $(1 + \alpha)$  when sampling the class mean vectors. When computing the NL score, the presumed between-class variances (those of the true x-vectors) are used, although the true data are sampled from a changed distribution. For comparison, we also compute the performance when the NL uses the changed between-class variances, which represents the performance when the domain mismatch is perfectly addressed (e.g., by retraining the NL parameters).

The results are shown in Figure 9. Comparing the difference between the red line (with domain mismatch) and the blue line (without domain mismatch), we can see if incorrect between-class variances are used, the NL performance is impacted, but not much.

In the next experiment, we simulate the case with an incorrect within-class variance. We multiply the original with-class variance by  $(1 + \alpha)$  when sampling the



data of each class for both enrollment and test. In computing the NL, the presumed within-class variance (that is 2.0 in our case) is used. Again, we compute the performance when the NL uses the changed within-class variance, which represents the performance when the domain mismatch is perfectly addressed.

The results are shown in Figure 10. Comparing the difference between the red line (with domain mismatch) and the blue line (without domain mismatch), we can see that when the within-class variance is incorrectly set, the performance is impacted, in particular when the true within-class variance is large but we assume it is small. The impact is more serious on the SV task compared to the SI task. This result suggests that a larger within-covariance is a safe choice when designing a practical system.



The final experiment simulates the shift on data, which is often observed when speaker recognition systems migrate to a new channel. We simply add a value  $\beta$  to all the dimensions of the sampled data, and then use the presumed NL parameters to compute the scores. The results are shown in Figure 11, where we also report the results without the shift. The results show that data shift impacts performance in a very significant way, and seems much more severe compared to the change on the between-class and within-class variances. The fatal impact of data shift has been reported with experiments on real datasets, e.g., [2].



#### 3.5.3 Domain adaptation

There are numerous studies on domain adaptation with PLDA (equivalent to NL based on a linear Gaussian model). The research can be categorized into three themes. The first theme adapts the covariances (or equivalently the factor loading matrices of PLDA) of the source domain to match the data in the target domain. This could be supervised or unsupervised. The supervised approach uses class labels in the target domain, and adapt the PLDA model following the Bayesian rule in principle [24, 61]. The unsupervised approach employs various clustering methods to generate pseudo classes, and then treat these pseudo classes as true speakers to conduct supervised adaptation [25]. The second theme analyzes the variation related to domains. This variation will be either removed from the data [1, 2, 30] or treated as a new subspace in LDA or PLDA models [47]. The third theme tries to learn a mapping function that transfers the data from the source domain to the target domain [53], or transfers data from multiple domains to a common domain [63].

Essentially, all these methods try to build a suitable statistical model for data in the target domain, by applying the knowledge of the source domain as much as possible, in the form of either model or data.

#### 3.6 Problem associated with enrollment-test condition mismatch

Another problem that may impact NL scoring in practice is the condition mismatch between enrollment and test. For example, one may enroll in an office but wants to perform test on the street. We will use simulation to investigate the impact of this enrollment-test condition mismatch.

Again, the simulation is based on the configuration of the x-vectors derived from VoxCeleb. For a better comparison, the within-class variance of the baseline (without any mismatch) is set to be 2.0. In each test, we sample 600 classes, and for each class, we sample one sample for enrollment and three samples for test. We run each test for 500 rounds, and report the averaged EER and IDR, plus the variations on them.

#### 3.6.1 Within-class variance mismatch

The simplest case is that the within-class variance changes during the test, but we compute the NL score using the within-class variance of the enrollment data. The results are shown in Figure 12, where the within-class variance of the test data is modified by multiplying the default value (i.e., the within-class variance of the enrollment data) by a scale factor  $1 + \alpha$ . We also report the performance using the (new) within-class variance of the test data for the NL scoring. Note that this is not a perfect solution as the new within-class variance matches the test data but does not match the enrollment data.

It can be seen that on both the SV and SI tasks, a larger within-class variance for the test data will lead to clear performance reduction, which is not surprising as a larger variance introduces more uncertainty. For SV, when the variance of the test data is larger than the enrollment variance, using the test variance (red curve) to compute the NL score leads to better performance compared to using the variance of the enrollment data. When the variance of the test data is smaller than that of the enrollment data, however, using the variance of the enrollment data (blue curve) seems slightly better. In other words, a larger within-class variance is preferred if there is a mismatch between the enrollment data and the test data. However, neither of these two choices is optimal: using the within-class variance of the enrollment data is not accurate for computing the prediction probability  $p(\boldsymbol{x}|\boldsymbol{\mu})$  of the test data and the normalization term  $p(\boldsymbol{x})$ , while using the within-class means, i.e.,  $p(\boldsymbol{\mu}|\boldsymbol{x}_1,...\boldsymbol{x}_{n_k})$ . We will present a condition transfer approach to solve this dilemma shortly. The new approach obtains the best performance, as shown by the brown curve in Figure 12.

For the SI task, we find that using the within-class variance of the enrollment data (blue curve) is better than using that of the test data (red curve). This is also expected, as in the SI task, the important thing is to estimate the class means, for which using the within-class variance of the enrollment data is theoretically correct. Once the class means are well estimated, using any within-class variance for test will lead to the same SI decision. In other words, the NL score is not impacted by the within-class variance mismatch on the SI task.

# 3.6.2 Mean scale and between-class variance mismatch

Another possible enrollment-test mismatch is that all the class means are scaled by a factor  $\alpha$  when generating the test data, while the within-class variance does not change. This scaling will lead to two problems: (1) *mean mismatch*: the class means of the enrollment data do not match the class means of the test data, leading



to incorrect likelihood  $p_k(\boldsymbol{x})$ ; (2) between-class distribution mismatch: the betweenclass variance of the test data is scaled in the same way as the class mean scaling. If we use the original between-class variance, the normalization term  $p(\boldsymbol{x})$  of the NL score will be inaccurate.

A simple compensation is to apply the same scale to the enrollment data. A problem with this compensation is that the scaling will change the within-class variance of the enrollment data. The ultimate effect will be the same as in the case of within-class variance mismatch: it would be a dilemma to choose the within-variance of the enrollment data or the test data.

Figure 13 shows the performance of five systems:

- Red curve: Scale the class means of the test data, and use the between-class and within-class variances of the enrollment data when computing the NL score. This is the result without any compensation.
- Blue curve: Scale the enrollment data and the class means of the test data in the same way, and use the between-class and within-class variances of the original enrollment data when computing the NL score. Since the enrollment data is scaled, the mean mismatch problem is solved, however the betweenclass distribution mismatch remains.
- Yellow curve: Scale the enrollment data and the class means of the test data in the same, and use the between-class variance of the *scaled* enrollment data (that is correct for both enrollment and test) and the within-class variance of the original enrollment data (that is correct for test but incorrect for enrollment) when computing the NL score. This approach solves the mean mismatch and between-class distribution mismatch, but the within-class variance is incorrect for enrollment.
- Purple curve: Scale the enrollment data and the class means of the test data in the same way, and use the between-class variance of the scaled enrollment data (that is correct for both enrollment and test) and the within-class variance of

the *scaled* enrollment data (that is correct for enrollment but incorrect for test) when computing the NL score. This approach solves the mean mismatch and between-class distribution mismatch, but the within-class variance is incorrect for test.

• Brown curve: Apply condition transfer that will be presented shortly.

From Figure 13, it can be seen that scaling the class means of the test data caused serious performance reduction, especially when the scale is large. Scaling the enrollment data to match the test data seems can mitigate this problem to a large extent. The impact of using the incorrect between-class and within-class distributions is not very substantial, indicating that mean mismatch is more serious compared to distribution mismatch in NL scoring. Finally, the condition transfer approach obtains the best performance, by correcting both the mean mismatch and the distribution mismatch.



# 3.6.3 Mean shift

In the third experiment, we shift the class means by  $\beta$  on each dimension when sampling the test data. This mean shift causes two problems for NL scoring: (1) *mean mismatch*: the class means of the enrollment data do not match the class means of the test data, leading to incorrect likelihood  $p_k(\boldsymbol{x}|\boldsymbol{\mu})$ , thus incorrect  $p_k(\boldsymbol{x})$ ; (2) between-class distribution shift: the between-class distribution of the test data is shifted in the same way as the mean shift, which leads to incorrect normalization  $p(\boldsymbol{x})$ .

A simple compensation is to shift the enrollment data in the same way as the test data. After the shift, the mean mismatch problem is mitigated, however the NL still uses the between-class distribution of the original enrollment data. This is essentially the data shift scenario in the domain mismatch experiment. Another compensation is to compute the posterior of the class mean  $p(\boldsymbol{\mu}|\boldsymbol{x})$  first, and then shift the mean of the posterior in the same way as the test data. By this way, the mean mismatch

is solved and the likelihood  $p_k(\boldsymbol{x})$  is correct, however the normalization  $p(\boldsymbol{x})$  is incorrect due to the shifted between-class distribution of the test data.

We report the results of four tests in Figure 14:

- Red curve: Shift the test data only. This is the case with no compensation.
- Blue curve: Shift the enrollment and test data in the same way. This is the data shift scenario in the domain mismatch experiment.
- Green curve: Shift the test data, and then shift the mean of the posterior  $p(\boldsymbol{\mu}|\boldsymbol{x})$ . It fully solves the mean mismatch problem, however the normalization  $p(\boldsymbol{x})$  is incorrect due to the shifted between-class distribution.
- Brown curve: No data shift.



The results shown in Figure 14 demonstrate that the mean shift on test data tends to cause significant performance degradation (red curve vs brown curve). This loss is comparable or even worse compared to the domain mismatch case (red curve vs blue curve). If we remove the mean mismatch but uses the incorrect normalization (green curve), the IDR performance recovers perfectly but the EER results become worse. The good performance on IDR is expected as the normalization term does not impact decisions of the SI task. The bad performance on EER demonstrates that an incorrect normalization may cause fatal performance loss on the SV task. An interesting observation is that for the EER results, removing the mean mismatch makes the performance even worse compared to doing nothing (green curve vs. red curve). This suggests that the errors caused by mean mismatch and between-class distribution shift are in opposite directions.

## 3.7 Condition transfer

We present a simple condition transfer approach based on the NL scoring, which is optimal under the linear Gaussian assumption. For simplicity, we will assume that the data have been shifted appropriately, so that the mean-shift problem does not exist. Denote the parameters of the NL models suitable for the enrollment and test data by  $\{\boldsymbol{\epsilon}, \sigma\}$  and  $\{\hat{\boldsymbol{\epsilon}}, \hat{\boldsymbol{\sigma}}\}$  respectively. Note that we have allowed a non-isotropic within-class covariance  $\mathbf{I}\hat{\boldsymbol{\sigma}}^2$  for the test data. Given enrollment samples  $\boldsymbol{x}_1^k, ... \boldsymbol{x}_{n_k}^k$  of class k, the posterior of its class mean will be computed using the between-class and within-class variances of the enrollment data:

$$p(\boldsymbol{\mu}_k | \boldsymbol{x}_1^k, \dots \boldsymbol{x}_{n_k}^k) = N(\boldsymbol{\mu}_k; \frac{n_k \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2} \bar{\boldsymbol{x}}_k, \mathbf{I} \frac{\boldsymbol{\epsilon}^2 \sigma^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2}).$$
(40)

Since there is no data shift, this posterior can be readily used to estimate the likelihood of the test sample, using the within-class variance  $\hat{\sigma}$  of the test data:

$$p_k(\boldsymbol{x}) = N(\boldsymbol{x}; \frac{n_k \boldsymbol{\epsilon}^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2} \bar{\boldsymbol{x}}_k, \mathbf{I}(\hat{\boldsymbol{\sigma}}^2 + \frac{\boldsymbol{\epsilon}^2 \sigma^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2})).$$
(41)

Augmented by the normalization term computed using the between-class and within-class variances of the test data, the NL score will have the following form:

$$\log NL(\boldsymbol{x}|k) \propto -||\frac{\boldsymbol{x} - \tilde{\boldsymbol{\mu}}_k}{\sqrt{\hat{\boldsymbol{\sigma}}^2 + \frac{\boldsymbol{\epsilon}^2 \sigma^2}{n_k \boldsymbol{\epsilon}^2 + \sigma^2}}}||^2 + ||\frac{\boldsymbol{x}}{\hat{\boldsymbol{\epsilon}}^2 + \hat{\boldsymbol{\sigma}}^2}||^2,$$
(42)

where

$$\tilde{\boldsymbol{\mu}}_{k} = \frac{n_{k} \boldsymbol{\epsilon}^{2}}{n_{k} \boldsymbol{\epsilon}^{2} + \sigma^{2}} \bar{\boldsymbol{x}}_{k}.$$
(43)

The condition transfer approach described above can be easily extended to handel more complex condition mismatch, which will be left for future work. Note that we have shown the performance of this method in Fig. 12 and Fig. 13. In both cases, it provides the best (actually optimal) performance.

# 4 Discussion

The NL formulation plays a central role in our simulation study. From the perspective of NL scoring, any performance loss can be attributed to data-model mismatch on the three components of the NL scoring: the enrollment model  $p(\boldsymbol{\mu}|\boldsymbol{x}_1,...,\boldsymbol{x}_{n_k})$ , the prediction model  $p(\boldsymbol{x}|\boldsymbol{\mu})$ , and the normalization model  $p(\boldsymbol{x})$ . The mismatch could be: (1) mismatch on distribution type (e.g., Gaussian assumed but Laplacian in reality); (2) mismatch on the mean (mean mismatch); (3) mismatch on the covariance (covariance mismatch). This analytical view provides a powerful and necessary tool for our simulation study. By this tool, we can analyze how a particular imperfection causes performance reduction, and design suitable algorithms to compensate for the impact, e.g., the conditional transfer algorithm.

The simulation results show that for a practical speaker recognition system, mean mismatch is the most risky. For example, in the data shift scenario of the domain mismatch experiment, the mean of the between-class distribution does not match the data, causing between-class mean mismatch; in the mean shift scenario of the enrollment-test condition mismatch experiment, the means of the within-class distributions of individual classes do not match the data, causing within-class mean mismatch. The performance reduction on these two scenarios is much more significant compared to on other scenarios.

Although our work focuses on the linear Gaussian NL, the NL formulation is general and can be easily extended by using nonlinear and non-Gaussian models, so that it deal with more complex data. Recently, we provide such an extension [38], by applying the invariance property of the NL score under invertible transforms, as discussed in Section 2.3. Specifically, we learn an invertible transform that maps the original data to a latent space where the data can be modeled by a linear Gaussian. According to the equivalence of the NL score in the original and the transformed space, this transform allows us using a linear Gaussian NL model to score data with a complex distribution. This is essentially a nonlinear extension of the PLDA model, which we call *neural discriminant analysis* (NDA). In our previous study, the NDA model produces very promising results [38].

The MBR optimum of the NL scoring may encourage more research on the speaker embedding approach. Since we have known that the NL score is MBR optimal, its performance will be ensured if the distribution of the speaker vectors meet the assumption of the model. This performance ensurance represents a clear advantage of the embedding approach compared to the so-called end-to-end approach [12, 27, 68]. Moreover, since the NL score is optimal if and only if the speaker vectors follow the assumed generative model, more research is encouraged on normalizing the speaker vectors, rather than pursuing other complicated scoring methods (e.g., discriminative PLDA [7]) or score calibration [15, 59]. Our recently work shows that speaker vector normalization is highly promising [9].

Finally, the main purpose of the paper is a full understanding for the NL score by simulation, so we have refrained from presenting any EER/IDR results on real SRE systems (large-scale experiments for the NL score with real data have been presented by other papers, e.g., [38]). We found that the simulation study is very useful and offers a lower bound and an upper bound for a potential technique. For the lower bound, it gives a clear justification that a technique does work if the presumed condition is matched, and so what we should do is to meet the condition. For the upper bound, it tells the maximum that a technique can achieve if the presumed condition is perfectly matched, so we should not intend to seek for more in real applications.

# 5 Conclusions

We present an analysis on the optimal score for speaker recognition based on the MAP principle and the linear Gaussian assumption. The analysis shows that the normalized likelihood (NL) is optimal for both identification and verification tasks in the sense of minimum Bayes risk. We also show that the NL score based on the linear Gaussian model is equivalent to the popular PLDA LR. The cosine score and Euclidean score can be regarded as two approximations of the optimal NL score. Comprehensive simulation experiments were conducted to study the behavior of the NL score, especially at the operation point of a true speaker recognition system. The major knowledge we obtained from the simulation study is that the NL performance may be seriously reduced by real-life imperfections, including the non-Gaussianality and non-homogeneity of the data, inaccurate estimation of the between-class and

within-class variances, and potential mismatch between enrollment and test conditions. Among all the detrimental factors, data shift caused the most significant performance reduction. We also proposed a condition transfer approach that can compensate for the enrollment-test mismatch.

# Abbreviations

BC: Between-class; DET: Detection error tradeoff; DNN: Deep neural net; EER: Equal error rate; GMM-UBM: Gaussian mixture model-universal background model; IDR: Identification rate; LDA: Linear discriminant analysis; LR: Likelihood ratio; MBR: Minimum Bayes risk; MFCC: Mel frequency cepstrum coefficient; NL: Normalized likelihood; PCA: Principle component analysis; PLDA: Probabilistic linear discriminant analysis; STD: Standard deviation; SI: Speaker identification; SV: Speaker verification; VAE: Variational auto-encoder; WC: Within-class.

# Declarations

# Availability of data and materials

The VoxCeleb dataset can be obtained from http://www.robots.ox.ac.uk/ ~vgg/data/voxceleb/. The CNCeleb dataset can be obtained from http://www.openslr.org/82/.

# Competing interests

The authors declare that they have no competing interests.

## Funding

This work was supported by the National Natural Science Foundation of China (NSFC) under the project No.61633013 and No.61371136.

# Authors' contributions

All the work was completed by DW.

#### Acknowledgements

Thanks to Dr. Lantian Li, Yunqi Cai and Zhiyuan Tang for the valuable discussion.

#### References

- Aronowitz H (2014) Compensating inter-dataset variability in PLDA hyper-parameters for robust speaker recognition. In: Proceedings of Odyssey: The Speaker and Language Recognition Workshop, pp 280–286
- Aronowitz H (2014) Inter dataset variability compensation for speaker recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 4002–4006
- Bai Z, Zhang XL, Chen J (2019) Partial AUC optimization based deep speaker embeddings with class-center learning for text-independent speaker verification. arXiv preprint arXiv:191108077
- 4. Bishop CM (2006) Pattern recognition and machine learning, Springer, chap 2
- 5. Blum A, Hopcroft J, Kannan R (2020) Foundations of data science. Cambridge University Press
- Borgström BJ, McCree A (2013) Discriminatively trained bayesian speaker comparison of i-vectors. In: 2013 IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, pp 7659–7662
- Burget L, Plchot O, Cumani S, Glembek O, Matějka P, Brümmer N (2011) Discriminatively trained probabilistic linear discriminant analysis for speaker verification. In: IEEE international conference on acoustics, speech and signal processing (ICASSP), IEEE, pp 4832–4835
- Cai W, Chen J, Li M (2018) Exploring the encoding layer and loss function in end-to-end speaker and language recognition system. In: Proceedings of Odyssey: The Speaker and Language Recognition Workshop, pp 74–81
- 9. Cai Y, Li L, Wang D, Abel A (2020) Deep normalization for speaker vectors. arXiv:200404095
- 10. Campbell JP (1997) Speaker recognition: A tutorial. Proceedings of the IEEE 85(9):1437–1462
- Chen N, Villalba J, Dehak N (2019) Tied mixture of factor analyzers layer to combine frame level representations in neural speaker embeddings. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 2948–2952

- rahman Chowdhury FR, Wang Q, Moreno IL, Wan L (2018) Attention-based models for text-dependent speaker verification. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 5359–5363
- 13. Chowdhury GG (2010) Introduction to modern information retrieval. Facet publishing
- 14. Chung JS, Nagrani A, Zisserman A (2018) VoxCeleb2: Deep speaker recognition. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1086–1090
- 15. Cumani S, Laface P (2019) Tied normal variance mean mixtures for linear score calibration. In: IEEE international conference on acoustics, speech and signal processing (ICASSP)
- Dehak N, Dehak R, Kenny P, Brümmer N, Ouellet P, Dumouchel P (2009) Support vector machines versus fast scoring in the low-dimensional total variability space for speaker verification. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH)
- Dehak N, Kenny PJ, Dehak R, Dumouchel P, Ouellet P (2011) Front-end factor analysis for speaker verification. IEEE Transactions on Audio, Speech, and Language Processing 19(4):788–798
- Ding W, He L (2018) MTGAN: Speaker verification through multitasking triplet generative adversarial networks. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 3633–3637
- 19. Dinh L, Krueger D, Bengio Y (2015) NICE: Non-linear independent components estimation. In: ICLR Workshop
- Dinh L, Sohl-Dickstein J, Bengio S (2016) Density estimation using real NVP. In: Neural Information Processing Systems - Deep Learning Symposium
- 21. Fan Y, Kang J, Li L, Li K, Chen H, Cheng S, Zhang P, Zhou Z, Cai Y, Wang D (2020) CN-CELEB: a challenging chinese speaker recognition dataset. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)
- Gao Z, Song Y, McLoughlin I, Li P, Jiang Y, Dai LR (2019) Improving Aggregation and Loss Function for Better Embedding Learning in End-to-End Speaker Verification System. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 361–365
- Garcia-Romero D, Espy-Wilson CY (2011) Analysis of i-vector length normalization in speaker recognition systems. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH)
- 24. Garcia-Romero D, McCree A (2014) Supervised domain adaptation for i-vector based speaker recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 4047–4051
- Garcia-Romero D, McCree A, Shum S, Vaquero C (2014) Unsupervised domain adaptation for i-vector speaker recognition. In: Proceedings of Odyssey: The Speaker and Language Recognition Workshop
- Hansen JH, Hasan T (2015) Speaker recognition by machines and humans: A tutorial review. IEEE Signal processing magazine 32(6):74–99
- 27. Heigold G, Moreno I, Bengio S, Shazeer N (2016) End-to-end text-dependent speaker verification. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 5115–5119
- Ioffe S (2006) Probabilistic linear discriminant analysis. In: European Conference on Computer Vision (ECCV), Springer, pp 531–542
- weon Jung J, Heo HS, ho Kim J, jin Shim H, Yu HJ (2019) RawNet: Advanced End-to-End Deep Neural Network Using Raw Waveforms for Text-Independent Speaker Verification. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1268–1272
- Kanagasundaram A, Dean D, Sridharan S (2015) Improving out-domain PLDA speaker verification using unsupervised inter-dataset variability compensation approach. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 4654–4658
- Kenny P (2010) Bayesian speaker verification with heavy-tailed priors. In: Proceedings of Odyssey: The Speaker and Language Recognition Workshop, p 14
- 32. Kingma DP, Dhariwal P (2018) Glow: Generative flow with invertible 1x1 convolutions. In: Advances in Neural Information Processing Systems (NIPS), pp 10,215–10,224
- Ko T, Peddinti V, Povey D, Seltzer ML, Khudanpur S (2017) A study on data augmentation of reverberant speech for robust speech recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 5220–5224
- Li L, Wang D, Xing C, Zheng TF (2016) Max-margin metric learning for speaker recognition. In: 10th International Symposium on Chinese Spoken Language Processing (ISCSLP), pp 1–4
- Li L, Chen Y, Shi Y, Tang Z, Wang D (2017) Deep speaker feature learning for text-independent speaker verification. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1542–1546
- Li L, Tang Z, Wang D, Zheng TF (2018) Full-info training for deep speaker feature learning. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 5369–5373
- Li L, Tang Z, Shi Y, Wang D (2019) Gaussian-constrained training for speaker verification. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 6036–6040
- Li L, Wang D, Zheng TF (2020) Neural discriminant analysis for speaker recognition. In: Submitted to Interspeech 2020
- Li R, Tuo NLD, Yu M, Su D, Yu D (2019) Boundary discriminative large margin cosine loss for text-independent speaker verification. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 6321–6325
- Lyu S, Simoncelli EP (2009) Nonlinear extraction of independent components of natural images using radial gaussianization. Neural computation 21(6):1485–1519
- 41. Mardia KV, Jupp PE (2009) Directional statistics, vol 494. John Wiley & Sons
- 42. McCree A, Sell G, Garcia-Romero D (2017) Extended variability modeling and unsupervised adaptation for plda speaker recognition. In: INTERSPEECH, pp 1552–1556
- Nagrani A, Chung JS, Zisserman A (2017) VoxCeleb: a large-scale speaker identification dataset. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH)

- Okabe K, Koshinaka T, Shinoda K (2018) Attentive statistics pooling for deep speaker embedding. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 2252–2256
- Povey D, Ghoshal A, Boulianne G, Burget L, Glembek O, Goel N, Hannemann M, Motlicek P, Qian Y, Schwarz P, et al (2011) The Kaldi speech recognition toolkit. In: IEEE workshop on automatic speech recognition and understanding
- Prince SJ, Elder JH (2007) Probabilistic linear discriminant analysis for inferences about identity. In: 2007 IEEE 11th International Conference on Computer Vision, IEEE, pp 1–8
- Rahman MH, Kanagasundaram A, Himawan I, Dean D, Sridharan S (2018) Improving PLDA speaker verification performance using domain mismatch compensation techniques. Computer Speech & Language 47:240–258
- Reynolds DA (2002) An overview of automatic speaker recognition technology. In: IEEE international conference on Acoustics, speech, and signal processing (ICASSP), pp 4072–4075
- Reynolds DA, Quatieri TF, Dunn RB (2000) Speaker verification using adapted Gaussian mixture models. Digital signal processing 10(1-3):19–41
- 50. Rudin W (2006) Real and complex analysis. Tata McGraw-hill education
- Sadjadi SO, Greenberg C, Singer E, Reynolds D, Mason L, Hernandez-Cordero J (2019) The 2018 NIST Speaker Recognition Evaluation. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1483–1487
- Salton G (1989) Automatic text processing: The transformation, analysis, and retrieval of. Reading: Addison-Wesley 169
- Shon S, Mun S, Kim W, Ko H (2017) Autoencoder based domain adaptation for speaker recognition under insufficient channel information. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1014–1018
- 54. Snyder D, Chen G, Povey D (2015) MUSAN: A Music, Speech, and Noise Corpus. ArXiv:1510.08484v1, 1510.08484
- Snyder D, Garcia-Romero D, Sell G, Povey D, Khudanpur S (2018) X-vectors: Robust DNN embeddings for speaker recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, pp 5329–5333
- Sra S (2018) Directional statistics in machine learning: a brief review. Applied Directional Statistics: Modern Methods and Case Studies p 225
- Stafylakis T, Rohdin J, Plchot O, Mizera P, Burget L (2019) Self-Supervised Speaker Embeddings. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 2863–2867
- Tu Y, Mak MW, Chien JT (2019) Variational Domain Adversarial Learning for Speaker Verification. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 4315–4319
- Van Leeuwen DA, Br N (2013) The distribution of calibrated likelihood-ratios in speaker recognition. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1619–1623
- Variani E, Lei X, McDermott E, Moreno IL, Gonzalez-Dominguez J (2014) Deep neural networks for small footprint text-dependent speaker verification. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 4052–4056
- Villalba J, Lleida E (2012) Bayesian adaptation of PLDA based speaker recognition to domains with scarce development data. In: Proceedings of Odyssey: The Speaker and Language Recognition Workshop, pp 47–54
- Wang J, Wang KC, Law MT, Rudzicz F, Brudno1 M (2019) Centroid-based deep metric learning for speaker recognition. In: International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 3652–3656
- Wang Q, Rao W, Sun S, Xie L, Chng ES, Li H (2018) Unsupervised domain adaptation via domain adversarial training for speaker recognition. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 4889–4893
- Wang S, Rohdin J, Burget L, Plchot O, Qian Y, Yu K, Cernocky J (2019) On the Usage of Phonetic Information for Text-Independent Speaker Embedding Extraction. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 1148–1152
- 65. Wang X, Li L, Wang D (2019) VAE-based domain adaptation for speaker verification. In: Proceedings of APSIPA ASC
- 66. Xie W, Nagrani A, Chung JS, Zisserman A (2019) Utterance-level aggregation for speaker recognition in the wild. In: IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp 5791–579
- Xing C, Wang D, Liu C, Lin Y (2015) Normalized word embedding and orthogonal transform for bilingual word translation. In: Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp 1006–1011
- Zhang SX, Chen Z, Zhao Y, Li J, Gong Y (2016) End-to-end attention based text-dependent speaker verification. In: Spoken Language Technology Workshop (SLT), IEEE, pp 171–178
- Zhang Y, Li L, Wang D (2019) VAE-based regularization for deep speaker embedding. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 4020–4024
- Zhou J, Jiang T, Li Z, Li L, Hong Q (2019) Deep Speaker Embedding Extraction with Channel-Wise Feature Responses and Additive Supervision Softmax Loss Function. In: Proceedings of the Annual Conference of International Speech Communication Association (INTERSPEECH), pp 2883–2887