

# A Unified Framework for Score Normalization Techniques Applied to Text Independent Speaker Verification

Johnny Mariéthoz and Samy Bengio

**Abstract**—The purpose of this paper is to unify several of the state-of-the-art score normalization techniques applied to text-independent speaker verification systems. We propose a new framework for this purpose. The two well-known Z- and T-normalization techniques can be easily interpreted in this framework as different ways to estimate score distributions. This is useful as it helps to understand the various assumptions behind these well-known score normalization techniques, and opens the door for yet more complex solutions. Finally, some experiments on the Switchboard database are performed in order to illustrate the validity of the new proposed framework.

**Index Terms**—speaker verification, score normalization, statistical framework, T-norm, Z-norm.

**EDICS Category: 2.SPEE**

## I. INTRODUCTION

Text-independent speaker verification systems have evolved through time [1]. The first systems had reasonable performance only in controlled conditions (no noise, same channel, same gender, etc). Over the years, researchers have improved their systems for unmatched conditions, thanks largely to score normalization techniques. In this paper, we propose a unified framework that explains several score normalization techniques used in text-independent speaker verification. Furthermore, an implementation of two of the most common techniques, the so-called T- and Z-normalization [2], is proposed here in this novel framework. While the two approaches are not strictly equivalent, in practice they give similar results. In fact, this new framework can be used to understand the assumptions that are implicit when using T- and Z-normalization. Moreover, it can also be used to develop new normalization techniques. The paper is organized as follows. In section II we present the classical framework used in speaker verification. In section III a new framework is proposed for score normalization. T- and Z-norm implementations in this framework are then given in sections IV-A and IV-B. Sections V and VIII show that the T- and Z-norm using this new framework are equivalent to their classical implementation. Finally we draw some conclusions.

## II. CLASSICAL FRAMEWORK USED IN SPEAKER VERIFICATION

Classical speaker verification models are based on a statistical framework. We are interested in  $P(S_i|X)$  the probability

that speaker  $S_i$  has pronounced sentence  $X$ . Using Bayes theorem, this can be expressed as follows:

$$P(S_i|X) = \frac{p(X|S_i)P(S_i)}{p(X)}. \quad (1)$$

In order to decide whether or not  $S_i$  has pronounced  $X$ , we compare  $P(S_i|X)$  to the probability that any other speaker has pronounced  $X$ , denoted  $P(\bar{S}_i|X)$ . When  $P(\bar{S}_i|X)$  is the same for all Clients, which is the assumption made in this paper, we replace it by a speaker independent model  $P(\Omega|X)$  where  $\Omega$  represents the *World* of all the speakers. The decision rule is then:

$$\text{if } P(S_i|X) > P(\Omega|X) \text{ then } X \text{ was uttered by } S_i. \quad (2)$$

Using equation (1), inequality (2) can be rewritten as:

$$\frac{p(X|S_i)}{p(X|\Omega)} > \frac{P(\Omega)}{P(S_i)} = \delta_i \quad (3)$$

where the ratio of the prior probabilities is usually replaced by a threshold  $\delta_i$  since it does not depend on  $X$  and is furthermore usually common for all speakers (hence  $\delta$ ). Taking the logarithm of (3) leads to the *log likelihood ratio* (LLR):

$$\text{llr}_i = \log \frac{p(X|S_i)}{p(X|\Omega)} > \log \delta_i = \Delta_i \approx \Delta. \quad (4)$$

## III. UNIFIED FRAMEWORK FOR SCORE NORMALIZATION

Most state-of-the-art text-independent speaker verification systems use linear score normalization functions of the form:

$$\text{llr}_{i_{norm}} = \frac{\text{llr}_i - \mu}{\sigma} > \Delta \quad (5)$$

where  $\mu$  and  $\sigma$  are respectively the mean and the standard deviation of a normal distribution of LLRs. These parameters are then estimated differently for each type of score normalizations. This paper proposes a unified framework for all kinds of normalization of the form of (5), and also other non-linear functions. We further propose an implementation for the two well-known T- and Z-normalization techniques.

We have seen that in text-independent speaker verification we are interested in the probability that a speaker  $S_i$  has pronounced a sentence  $X$ . Let us now consider the LLR as an additional random variable, and let us introduce it in the original framework by looking at  $P(S_i|X, \text{llr}_i)$ , the probability that a speaker  $S_i$  has pronounced a sentence  $X$  and obtained

an LLR of  $\text{llr}_i$ . Using the same approach as in section II, we obtain:

$$P(S_i|\text{llr}_i, \mathbf{X}) > P(\Omega|\text{llr}_i, \mathbf{X}). \quad (6)$$

Using inequality (6) and the Bayes theorem, it can then be rewritten as:

$$p(\text{llr}_i, \mathbf{X}|S_i) \frac{P(S_i)}{p(\text{llr}_i, \mathbf{X})} > p(\text{llr}_i, \mathbf{X}|\Omega) \frac{P(\Omega)}{p(\text{llr}_i, \mathbf{X})}. \quad (7)$$

Applying some simplifications to inequality (7) yields:

$$\frac{p(\text{llr}_i, \mathbf{X}|S_i)}{p(\text{llr}_i, \mathbf{X}|\Omega)} > \frac{P(\Omega)}{P(S_i)}. \quad (8)$$

Using inequality (8) and the conditional law of probabilities gives:

$$\frac{p(\text{llr}_i|\mathbf{X}, S_i) p(\mathbf{X}|S_i)}{p(\text{llr}_i|\mathbf{X}, \Omega) p(\mathbf{X}|\Omega)} > \frac{P(\Omega)}{P(S_i)}. \quad (9)$$

Taking the logarithm of inequality (9), we finally obtain:

$$\text{llr}'_i = \log \frac{p(\text{llr}_i|\mathbf{X}, S_i)}{p(\text{llr}_i|\mathbf{X}, \Omega)} + \text{llr}_i > \log \frac{P(\Omega)}{P(S_i)} \approx \Delta. \quad (10)$$

Comparing equation (10) of this new framework with the original equation (4), we can see that a new term appears. It is the log of the ratio of two likelihoods estimated by two score distributions. The numerator represents the distribution of LLRs for a given access  $\mathbf{X}$  and for client  $S_i$ . The denominator represents the distribution of LLRs for a given access  $\mathbf{X}$  and for all impostors  $\Omega$ . We will see that, depending on how these two distributions are estimated, we can obtain classical score normalization techniques such as T-norm (when estimated on a test access) or Z-norm (when estimated for each client  $S_i$ ).

#### IV. RELATION TO EXISTING NORMALIZATION TECHNIQUES

##### A. T-norm

The T-norm, as introduced in [2] and [3], estimates  $\mu$  and  $\sigma$  as the mean and the standard deviation of the log likelihood ratios (LLRs) using models of a subset of impostors, for a particular test access  $\mathbf{X}_0$ .

$$\mu_N = \frac{1}{N} \sum_n \text{llr}_n(\mathbf{X}_0) \quad (11)$$

$$\sigma_N = \sqrt{\frac{1}{N} \sum_n (\text{llr}_n(\mathbf{X}_0) - \mu_N)^2} \quad (12)$$

where  $N$  is the number of impostor models and  $\text{llr}_n$  is the score for the  $n^{\text{th}}$  impostor model for the particular access  $\mathbf{X}_0$ . Using (5) we obtain:

$$\text{llr}_{it-norm} = \frac{\text{llr}_i - \mu_N}{\sigma_N} > \Delta. \quad (13)$$

Let us now show how it is possible to perform T-normalization using our new framework under reasonable assumptions. We also show in the Appendix a comparison of our framework and the T-norm implementation found in the literature.

Given the framework described in section III, we must define two distributions, which will be here defined as Normal, as follows:

$$\hat{p}(\text{llr}|\mathbf{X}, S_i) = \mathcal{N}(\text{llr}_{S_i}; \mu_{S_i}, \sigma_{S_i}) \quad (14)$$

$$\hat{p}(\text{llr}|\mathbf{X}, \Omega) = \mathcal{N}(\text{llr}_{S_i}; \mu_{\Omega}, \sigma_{\Omega}) \quad (15)$$

where  $\mu_{S_i}, \sigma_{S_i}$  are the parameters of the client distribution and  $\mu_{\Omega}, \sigma_{\Omega}$  are the parameters of the impostor distribution. To obtain the T-norm we make the assumption that the standard deviations are equal:  $\sigma_N = \sigma_{S_i} = \sigma_{\Omega}$ . We thus obtain:

$$\begin{aligned} \log \frac{\hat{p}(\text{llr}|\mathbf{X}, S_i)}{\hat{p}(\text{llr}|\mathbf{X}, \Omega)} &= -\frac{1}{2\sigma_N^2} \left( (\text{llr}_{S_i} - \mu_{S_i})^2 - (\text{llr}_{S_i} - \mu_{\Omega})^2 \right) \\ &\quad - \log \frac{\sqrt{2\pi\sigma_N^2}}{\sqrt{2\pi\sigma_N^2}} \\ &= \frac{\mu_{S_i} - \mu_{\Omega}}{\sigma_N^2} \left( \text{llr}_{S_i} - \frac{\mu_{S_i} + \mu_{\Omega}}{2} \right). \end{aligned} \quad (16)$$

If we now define the means as:

$$\begin{aligned} \mu_{S_i} &= \text{llr}_{S_i} \\ \mu_{\Omega} &= \mu_N \end{aligned} \quad (17)$$

we obtain

$$\text{llr}_{S_i} + \frac{(\text{llr}_{S_i} - \mu_N)^2}{2\sigma_N^2} > \Delta. \quad (18)$$

Note that equations (17) and (18) are valid only when  $\text{llr}_{S_i} > \mu_N$ . A reasonable thing to do is to reject directly without any normalization a claimed speaker if its obtained LLR is smaller than the average of LLRs over a subset of impostors. The consequence of this on the T-norm equation is to force the threshold  $\Delta$  in (13) to be positive.

##### B. Z-norm

The basis of Z-norm [2] is to test a speaker model against example impostor utterances and to use the corresponding LLR scores to estimate a speaker specific mean and standard deviation:

$$\mu_J = \frac{1}{J} \sum_j \text{llr}_{S_i}(\mathbf{X}_j) \quad (19)$$

$$\sigma_J = \sqrt{\frac{1}{J} \sum_j (\text{llr}_{S_i}(\mathbf{X}_j) - \mu_J)^2} \quad (20)$$

where  $J$  is the number of impostor accesses. Using a similar approach to that in section IV-A, the estimate of the two distributions needed for the proposed unified framework becomes:

$$\hat{p}(\text{llr}|\mathbf{X}, S_i) = \mathcal{N}(\text{llr}_{S_i}; \mu_{S_i}, \sigma_{S_i}) \quad (21)$$

$$\hat{p}(\text{llr}|\mathbf{X}, \Omega) = \mathcal{N}(\text{llr}_{S_i}; \mu_{\Omega}, \sigma_{\Omega}) \quad (22)$$

with, again, the same standard deviation,  $\sigma_J = \sigma_{S_i} = \sigma_{\Omega}$ . If we now define the means as follows:

$$\begin{aligned} \mu_{S_i} &= \text{llr}_{S_i} \\ \mu_{\Omega} &= \mu_J \end{aligned} \quad (23)$$

then using equations (23) and (16) we obtain:

$$\text{llr}_{S_i} + \frac{(\text{llr}_{S_i} - \mu_J)^2}{2\sigma_J^2} > \Delta. \quad (24)$$

Finally, as explained for the T-norm at the end of section IV-A, we also need to reject a claimed access if  $(\|r_{s_i} < \mu_J)$ .

## V. EXPERIMENTS

The goal of these experiments is to show that the proposed framework can indeed be used to perform T-norm or Z-norm while obtaining the same performance as the original methods, and, gaining some insight about the underlying assumptions.

### A. Performance Measure in Speaker Verification

The performance of a speaker verification system is usually represented in terms of false acceptance rate (FAR, the number of false acceptances divided by the number of impostor accesses) and false rejection rate (FRR, the number of false rejections divided by the number of client accesses). A summary of these two values is often given by the half total error rate (HTER, the average of FAR and FRR), or by the equal error rate (EER, the point where FAR is equal to FRR<sup>1</sup>). It is also possible to represent graphically the performance using DET curves [4] which, similarly to ROC curves, show FAR with respect to FRR for various values of the threshold of equation (4), but with a Normal scale transformation. More recently, a new Expected Performance Curve (EPC) has been proposed by [5], which has shown to provide a fairer comparison between models. The procedure optimizes a convex combination of the individual performance measures,  $ep = \alpha FAR + (1 - \alpha) FRR$ , for various values of  $\alpha \in [0, 1]$  on the validation set during the training procedure used for parameter selection, and then plots HTER on the test set as a function of  $\alpha$ . In this way, each point on the graph contains its underlying *a priori* threshold selection procedure and is thus comparable to similar points (same  $\alpha$ ) coming from other models. Thus, this curve can be seen as an *a priori* DET curve.

### B. Database and Protocol

The comparison was done on a subset of the database that was used for the *NIST 2000 Speaker Recognition Evaluation*, which comes from the Switchboard-2 Phase 1 and 2 Corpus collected by the Linguistic Data Consortium. This data was used as an *evaluation* set while the World model and the *development* data come from previous NIST campaigns. While in the original database two different handsets were used (carbon and electret), in the subset selected for this paper, we only used data from electret handsets. This protocol was first proposed by the ELISA consortium as a reference for the NIST 2001 evaluation. We separated the data into male and female data, in order to create two different World models. The male World model was trained on 137 speakers for a total of 1.5 hours of speech, while the female World model was trained on 218 speakers for a total of 3 hours of speech. After that, the two World models were merged: the new World model has the same mean and variance vectors as the concatenation of the two gender dependent World models and the weights

<sup>1</sup>Note that EER is an *a posteriori* measure in the sense that the underlying threshold is necessarily chosen on the test set.

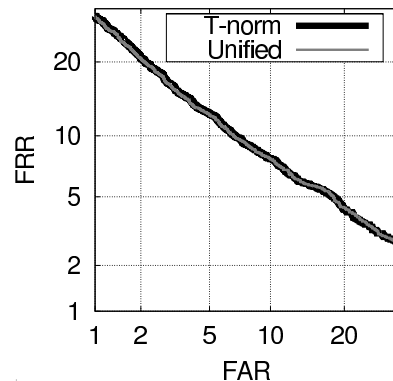


Fig. 1. DET curves on the NIST 2000 evaluation set for the T-norm and unified framework T-norm systems.

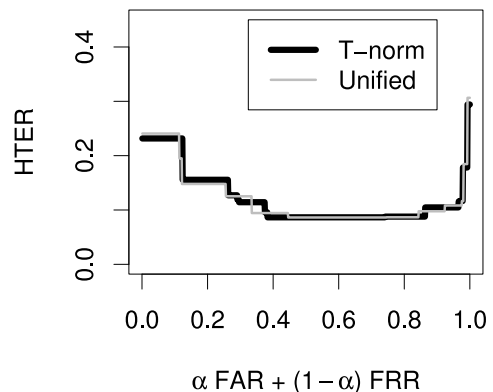


Fig. 2. EPC curves on the NIST 2000 evaluation set for the T-norm and unified framework T-norm systems.

are normalized in order to satisfy the constraint that they should sum to 1. For both development and evaluation Clients, approximately 2 minutes of telephone speech were used to train the models and each test access was less than 1 minute long. The development population consisted of 45 males, with 417 males in the evaluation set. The total number of accesses in the development population was 2441 and 27893 for the evaluation population with a proportion of 10% of true target accesses.

### C. Experimental Results

To verify the validity of our framework and the underlying assumptions, we first compared the standard Z- and T-normalizations and the version derived from the proposed framework. Figure 1 and 3 present the results using DET curves since these curves are often used in the literature. Unfortunately, as explained in section V-A, DET curves do not take into account the threshold estimation procedure. We thus also present results using EPC in Figure 2 and 4. In both cases the two curves match each other. These results show that the two approaches are equivalent.<sup>2</sup>

<sup>2</sup>In fact they are perfectly equal if we remove  $\|r_{s_i}$  in equation (18) and (24).

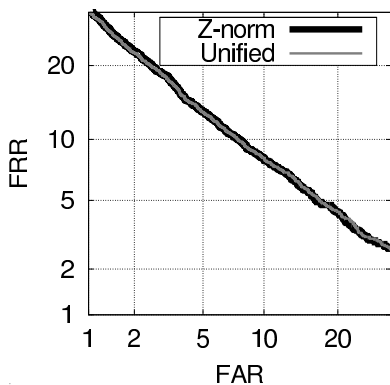


Fig. 3. DET curves on the NIST 2000 evaluation set for the Z-norm and unified framework Z-norm systems.

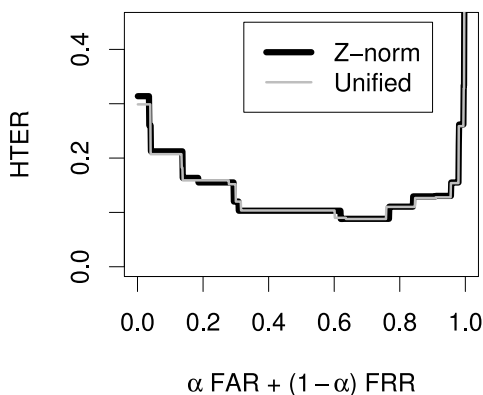


Fig. 4. EPC curves on the NIST 2000 evaluation set for the Z-norm and unified framework Z-norm systems.

## VI. CONCLUSION

In this paper, we have presented a new unified framework for text-independent speaker verification score normalization techniques. We have shown that the T- and Z-normalizations can be formalized using this new framework. Theoretical and empirical results show that the implementation found in the literature for T- and Z-norm are equivalent to our implementation. This helps to interpret T- and Z-norm as a way to estimate score distributions using two Normal distributions with the same variance. These normalization techniques have a very simple form in this framework and we can thus hope to find an even better estimate of LLR distributions. Indeed, there is no reason to force the LLR distribution to be Normally distributed, as done for the T- and Z-norm. Using our framework it is possible to approximate these distributions using more complex models, such as Mixtures of Gaussians for example. We hope that this framework will be used to propose new score normalization methods and also to improve understanding of this type of algorithms.

## VII. ACKNOWLEDGMENTS

This research has been partially carried out in the framework of the Swiss NCCR project (IM)2. It was also supported in part by the IST Program of the European Community, under the PASCAL Network of Excellence, IST-2002-506778, funded

in part by the Swiss OFES. All experiments were performed using the *Torch* package [6]. We would also like to thank Iain McCowan and David Granger for fruitful discussion and proofreading.

## VIII. APPENDIX

In this Appendix, we show the difference between the T-norm implementation found in the literature and our implementation using a unified framework. This demonstration can also be applied to Z-normalization.

The new implementation is given by:

$$\text{llr}_{S_i} + \frac{(\text{llr}_{S_i} - \mu_N)^2}{2\sigma_N^2} > \Delta \quad (25)$$

The classical method to implement T-norm is equivalent to the second term of the left side of equation (25) since:

$$\begin{aligned} \frac{(\text{llr}_{S_i} - \mu_N)^2}{2\sigma_N^2} &> \Theta \\ (\text{llr}_{S_i} - \mu_N)^2 &> \Theta * 2\sigma_N^2 \\ (\text{llr}_{S_i} - \mu_N)^2 - 2\Theta * \sigma_N^2 &> 0 \\ \left[ (\text{llr}_{S_i} - \mu_N - \sqrt{2\Theta} * \sigma_N) \right. \\ &\cdot \left. (\text{llr}_{S_i} - \mu_N + \sqrt{2\Theta} * \sigma_N) \right] > 0 \end{aligned} \quad (26)$$

and if  $\text{llr}_{S_i} > \mu_N$  then we can simplify (26) further into:

$$\begin{aligned} \text{llr}_{S_i} - \mu_N - \sqrt{2\Theta} * \sigma_N &> 0 \\ \frac{\text{llr}_{S_i} - \mu_N}{\sigma_N} &> \sqrt{2\Theta}. \end{aligned} \quad (27)$$

This inequation has a real solution only when  $\Theta > 0$ , which is true if  $\text{llr}_{S_i} > \mu_N$ . This assumption is reasonable: we do not want to accept an access if the LLR on the client model is smaller than the average LLR obtained over a subset of impostors. Given this reasonable assumption we can see the standard T-norm as a simplification of the T-norm using our new unified framework.

## REFERENCES

- [1] F. Bimbot, J.-F. Bonastre, C. Fredouille, G. Gravier, I. Magrin-Chagnolleau, S. Meignier, T. Merlin, J. Ortega-Garcia, D. Petrovsk-Delacrtaz, and D. Reynolds, "A tutorial on text-independent speaker verification," *EURASIP Journal on Applied Signal Processing*, vol. 4, pp. 430–451, 2004.
- [2] R. Auckenthaler, M. Carey, and H. Lloyd-Thomas, "Score normalization for text-independent speaker verification systems," *Digital Signal Processing*, vol. 10, pp. 42–54, January/April/July 2000.
- [3] J. Navratil and G. N. Ramaswamy, "The awe and mystery of t-norm," in *Proc. of the European Conference on Speech Communication and Technology*, 2003, pp. 2009–2012.
- [4] A. Martin, G. Doddington, T. Kamm, M. Ordowski, and M. Przybocki, "The DET curve in assessment of detection task performance," in *Proceedings of Eurospeech'97, Rhodes, Greece*, 1997, pp. 1895–1898.
- [5] S. Bengio and J. Mariéthoz, "The expected performance curve: a new assessment measure for person authentication," in *Proceedings of Odyssey 2004: The Speaker and Language Recognition Workshop*, 2004.
- [6] R. Collobert, S. Bengio, and J. Mariéthoz, "Torch: a modular machine learning software library," IDIAP, Technical Report IDIAP-RR 02-46, 2002.