

Discriminative Subspace Modeling of SNR and Duration Variabilities for Robust Speaker Verification

Na LI^a, Man-Wai MAK^a, Wei-Wei LIN^a, Jen-Tzung CHIEN^b

^a *Dept. of Electronic and Information Engineering, The Hong Kong Polytechnic University, Hong Kong SAR of China*

^b *Dept. of Electrical and Computer Engineering, National Chiao Tung University, Taiwan*

Abstract

Although i-vectors together with probabilistic LDA (PLDA) have achieved a great success in speaker verification, how to suppress the undesirable effects caused by the variability in utterance length and background noise level is still a challenge. This paper aims to improve the robustness of i-vector based speaker verification systems by compensating for the utterance-length variability and noise-level variability. Inspired by the recent findings that noise-level variability can be modeled by a signal-to-noise ratio (SNR) subspace and that duration variability can be modeled as additive noise in the i-vector space, we propose to add an SNR factor and a duration factor to the PLDA model. In this framework, we assume that i-vectors derived from utterances with comparable durations share similar duration-specific information and that i-vectors extracted from utterances within a narrow SNR range have similar SNR-specific information. Based on these assumptions, an i-vector can be represented as a linear combination of four components: speaker, SNR, duration, and channel. A variational Bayes algorithm is developed to infer this latent variable model via a discriminative subspace training procedure. In the testing stage, different variabilities are compensated when computing the likelihood ratio. Experiments on Common Conditions 1 and 4 in NIST 2012 SRE show that the proposed model outperforms the conventional PLDA and SNR-invariant PLDA. Results also show that the proposed model performs better than the uncertainty-propagation PLDA (UP-PLDA) for long test utterances.

Keywords: Speaker verification, duration variation, SNR mismatch, variational Bayes, I-vector, PLDA

1 Introduction

In text-independent speaker verification, i-vectors [6] have become the most popular feature representation in recent years. Inspired by the joint factor analysis (JFA) [15, 17, 18] framework, both the speaker and undesirable information (e.g. channel, additive noise, and so on) were compressed into a low-dimensional subspace called the total variability subspace, through which utterances with variable durations can be represented as low-dimensional i-vectors of fixed-length. Such a representation converts a speaker verification problem to an ordinary biometric pattern recognition problem similar to face recognition and fingerprint recognition. Based on the i-vector representation, many statistical techniques have been applied to deal with the mismatch between the training and test utterances. For example, linear discriminant analysis (LDA) [2] followed by within-class covariance normalization (WCCN) [12] were applied to i-vectors to compensate for session variability; then cosine distance between the target speaker's i-vector and test i-vector was used as the similarity measure between the target speaker and the test speaker. More recently, probabilistic LDA (PLDA) [31] was employed to suppress the channel- and session-variability within the i-vector space. Typically, i-vectors were preprocessed by a series of transformations – WCCN, length normalization [7], and LDA – before presenting the i-vectors to a Gaussian PLDA model.

Although the i-vector/PLDA framework performs well in suppressing session variability, it still has the following limitations: (1) the ability of PLDA in modeling the variability arising from utterances of different SNRs is limited; (2) i-vectors extracted from short utterances are less reliable than those extracted from long utterances [19], leading to performance degradation when only short utterances are available.

One of the main focuses of NIST 2012 SRE is robust speaker verification in which the SNR and length of enrollment and test utterances have substantial variation. To improve the noise robustness of i-vector/PLDA systems, several methods have been proposed. In [10], clean and noisy utterances were pooled together to train a robust PLDA model. Garcia-Romero *et al.* [8] employed multi-condition training to train multiple PLDA models, one for each condition. A robust system was then constructed by combining all of the PLDA models according to the posterior probability of each condition. In [23, 24], a mixture of SNR-dependent PLDA was proposed so that each mixture focuses on a small range of SNRs. During verification, the mixtures cooperated with each other to deal with utterances of various noise levels. By assuming that i-vectors derived from utterances falling within a narrow

40 SNR range should share similar SNR-specific information, we have recently
41 proposed to add an SNR-subspace to the conventional PLDA model, result-
42 ing in SNR-invariant PLDA [21, 20]. With the added SNR subspace, the
43 SNR-invariant PLDA can capture both speaker, noise-level, and channel
44 variabilities embedded in the i-vectors.

45 The problem of duration variability in utterances has attracted atten-
46 tion in the community because an i-vector extracted from a short utterance
47 should not be treated as being equally reliable as an i-vector extracted from
48 a long utterance. The reason is that the posterior distribution of hidden
49 variables in the i-vector extractor is a Gaussian whose covariance matrix is
50 related to the utterance duration. The shorter the utterance is, the larger
51 the covariance will become, leading to greater uncertainty in the estimated
52 i-vector.

53 The issue of duration variability has been addressed to a certain ex-
54 tent in the past. For example, Sarkar *et al.* [32] investigated how duration
55 mismatches affect the optimal choice of the duration of training utterances
56 for estimating the parameters of i-vector systems. In [19], the uncertainty
57 arising from the i-vector extraction process was propagated into a PLDA
58 model. This method did not treat an i-vector as the maximum *a posteri-*
59 *ori* point estimate, but rather as a random vector whose uncertainty was
60 represented by the posterior covariance matrix of the latent factors. The
61 shorter the utterance, the larger the posterior covariances. By propagating
62 this information into PLDA and using a loading matrix to model the vari-
63 ability due to duration variation, the resulting PLDA model better handled
64 the length-variability than the conventional PLDA model. Cumani *et al.*
65 [5, 4] did not map an utterance to a single i-vector, but instead mapped it
66 to the posterior distribution of i-vectors. Then, the likelihood of two speech
67 segments coming from the same speaker was obtained by integrating out all
68 possible i-vectors based on the i-vector posterior density.

69 Hasan *et al.* [11] found that duration variability could be modeled as ad-
70 ditive noise in the i-vector space. A short-utterance variance normalization
71 technique and a short-utterance variance modeling approach were proposed
72 in [14] to compensate for utterance-length variability. In [34], a weight
73 associated with the utterance’s duration was added to the corresponding
74 i-vector; then duration-weighted means, covariance matrix, and within-class
75 scatter matrix were computed; finally, principal component analysis (PCA)
76 and WCCN were applied using these duration-weighted terms to take utter-
77 ance duration into account. Motivated by the belief that i-vectors derived
78 from long utterances are more reliable [19] and therefore their corresponding
79 covariances in the PLDA model should be smaller, Cai *et al.* [3] proposed to

80 regularize the PLDA covariance matrix by scaling it by a duration-dependent
81 exponential term. On top of this duration-dependent covariance regular-
82 ization, Hong et al. [13] introduced a quality measure function for score
83 calibration, which effectively compensated for the score shift due to dura-
84 tion mismatch. In [36], a denoising autoencoder was used to compensate
85 for the phonetic imbalance in short utterances. Given a short utterance, the
86 autoencoder received an i-vector and a phonetic vector (the utterance’s zero-
87 order statistics) as input and produced an output comprising an i-vector as
88 if it were produced by a phonetically balance utterance. The autoencoder
89 was trained by using the i-vectors and phonetic vectors derived from many
90 short-long utterance pairs.

91 This paper focuses on improving the robustness of the state-of-the-art
92 i-vector/PLDA systems when duration mismatch and SNR mismatch be-
93 tween the training and test utterances occur simultaneously. According to
94 [11, 14], duration variability in the i-vectors can be modeled as additive
95 noise in i-vector space. If the i-vector extracted from a long utterance is
96 considered as “clean”, the i-vector extracted from a short utterance can
97 be considered as “noisy”. Inspired by this observation, we propose a new
98 method to deal with the mismatch caused by the variabilities in SNR and
99 duration. Our proposal is motivated by the success of SNR-invariant PLDA
100 in dealing with SNR mismatch [21, 20]. More specifically, we attempt to
101 make the i-vector/PLDA framework more resilient to SNR and duration
102 variabilities by introducing two discriminant subspaces – namely SNR sub-
103 space and duration subspace – to the PLDA models. These subspaces are
104 trained discriminatively by exploiting the SNR and duration information in
105 the training utterances. Through joint discriminative training, these sub-
106 spaces enable the new PLDA models to capture not only speaker and channel
107 variabilities, but also SNR and duration variabilities. In the proposed model,
108 the speaker component, SNR component, and duration component live in
109 three different subspaces which can be inferred according to the variational
110 Bayes procedure. During the verification stage, SNR variability, duration
111 variability, and channel variability are marginalized out when the likelihood
112 ratio is computed.

113 The organization of this paper is as follows: Section 2 describes the i-
114 vector/PLDA speaker verification. Based on different assumptions, a new
115 method of estimating the parameters of duration-invariant PLDA and two
116 new scoring methods are proposed in Section 3. The proposed modeling
117 method, namely SNR- and duration-invariant PLDA, is explained for robust
118 speaker verification in Section 4. The experimental results and analysis of
119 the proposed framework are detailed in Section 5 and Section 6, respectively.

120 Finally, conclusions are drawn in Section 7.

121 **2. I-vector/PLDA Speaker Verification**

122 *2.1. Conventional PLDA*

123 In the conventional i-vector/PLDA framework [16], an i-vector \mathbf{x}_{ij} is
 124 regarded as an observation generated from a linear model [31, 30]:

$$\mathbf{x}_{ij} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{G}\mathbf{r}_{ij} + \boldsymbol{\epsilon}_{ij} \quad (1)$$

125 where \mathbf{m} is the global mean of i-vectors, \mathbf{V} defines the speaker subspace, \mathbf{G}
 126 defines the channel subspace, \mathbf{h}_i and \mathbf{r}_{ij} are the latent factors depending on
 127 the speaker and session respectively, and $\boldsymbol{\epsilon}_{ij}$ denotes a residual term which
 128 follows a Gaussian distribution $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \boldsymbol{\Sigma})$. Typically, $\boldsymbol{\Sigma}$ is a diagonal matrix
 129 aiming to model any remaining variation that cannot be described by $\mathbf{V}\mathbf{V}^\top$
 130 and $\mathbf{G}\mathbf{G}^\top$.

131 According to [16, 7], the PLDA model in Eq. 1 can be divided into two
 132 parts: (1) the speaker part ($\mathbf{m} + \mathbf{V}\mathbf{h}_i$) that depends on the i -th speaker
 133 only and (2) the channel part ($\mathbf{G}\mathbf{r}_{ij} + \boldsymbol{\epsilon}_{ij}$) that depends not only on the
 134 i -th speaker but also on the j -th session. As i-vectors are of sufficiently
 135 low dimension, the term $\mathbf{G}\mathbf{r}_{ij}$ can be absorbed into $\boldsymbol{\Sigma}$ if the latter is a full
 136 covariance matrix. Accordingly, the Gaussian PLDA model can be simplified
 137 as follows [33]:

$$\mathbf{x}_{ij} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \boldsymbol{\epsilon}_{ij}, \quad (2)$$

138 where $\boldsymbol{\epsilon}_{ij} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma}$ being a full covariance matrix. This paper
 139 adopts this simplified model.

140 *2.2. SNR-invariant PLDA*

141 To enhance the robustness of i-vector/PLDA, we have recently proposed
 142 an SNR-invariant PLDA model (SI-PLDA) [21, 20] to deal with SNR mis-
 143 match. In this model, training utterances are first divided into K groups
 144 according to their SNRs. As a result, each of the training i-vectors is asso-
 145 ciated with one SNR group. Denote \mathbf{x}_{ij}^k as the j -th i-vector from speaker i
 146 in the k -th SNR group. Then, \mathbf{x}_{ij}^k is expressed as:

$$\mathbf{x}_{ij}^k = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{U}\mathbf{w}_k + \boldsymbol{\epsilon}_{ij}^k, \quad (3)$$

147 where \mathbf{m} is the global mean of i-vectors, \mathbf{V} defines the speaker subspace, \mathbf{h}_i
 148 is a latent speaker factor with a standard normal prior, \mathbf{U} defines the SNR
 149 subspace, \mathbf{w}_k is a latent SNR factor with a standard normal prior, $\boldsymbol{\epsilon}_{ij}^k$ is a

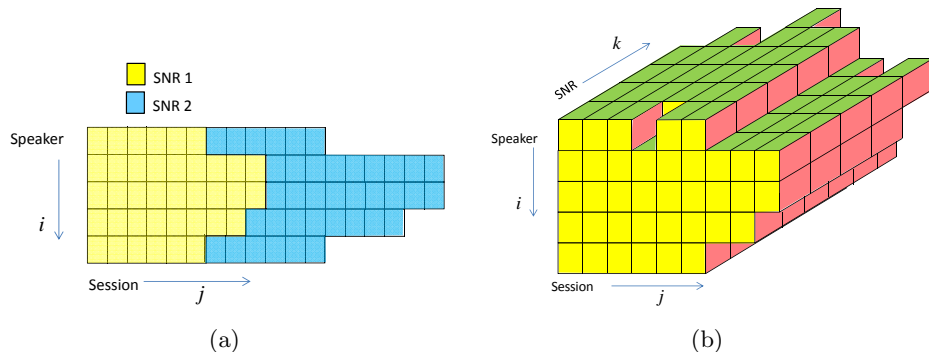


Figure 1: (a) Arrangement of training i-vectors in the multi-condition training of conventional PLDA. Each small square represents an i-vector. While the training set comprises two SNR groups, PLDA training ignores the group labels and sums over the statistics across both groups. (b) Arrangement of training i-vectors in SNR-invariant PLDA. Each small cube represents an i-vector. For the i -th speaker, there are $H_i(k)$ i-vectors from the k -th SNR group. Training in SNR-invariant PLDA considers the group labels and sums over the statistics within individual groups.

150 residual term with distribution $\mathcal{N}(\epsilon|\mathbf{0}, \mathbf{\Sigma})$. In [21, 20], $\mathbf{\Sigma}$ is a full covariance
 151 matrix aiming to model the channel variability.

152 The key difference between the conventional PLDA (Eq. 1) and SNR-
 153 invariant PLDA (Eq. 3) is that the former uses a channel subspace (\mathbf{G}) to
 154 model channel variability, whereas the latter uses an SNR subspace (\mathbf{U}) to
 155 capture the variability due to noise level differences. As a result, the SNR
 156 latent factors (\mathbf{w}_k in Eq. 3) depend on the SNR groups, whereas the session
 157 latent factor (\mathbf{r}_{ij} in Eq. 1) depends on the speaker and session.

158 Fig. 1 illustrates how the labels (speaker and SNR groups) can be used
 159 in training these two types of PLDA models. As can be seen, in the conven-
 160 tional PLDA (Fig. 1(a), Eq. 1, and Eq. 2), the i-vectors for each speakers
 161 are treated equally regardless of which SNR group they come from. On
 162 the other hand, in SI-PLDA (Fig. 1(b) and Eq. 3), i-vectors derived from
 163 utterances of similar SNR are grouped together in a vertical slice. These
 164 extra SNR labels, together with the speaker labels, help to suppress the
 165 SNR variability in the i-vectors.

166 3. Duration-invariant PLDA

167 According to [11, 14], duration variability in the i-vectors can be modeled
 168 as additive noise in i-vector space. Inspired by the success of SI-PLDA in

169 handling SNR variability, we propose to handle duration variability by a
 170 duration-invariant PLDA (DI-PLDA).

171 3.1. Generative Model and EM Formulation

172 Assume that we have a set of i-vectors

$$\mathcal{X} = \{\mathbf{x}_{ij}^p | i = 1, \dots, S; j = 1, \dots, H_i(p); p = 1, \dots, P\}$$

173 obtained from S speakers, where \mathbf{x}_{ij}^p is the j -th utterance from speaker i
 174 at the p -th duration group. For the i -th speaker, there are $H_i(p)$ i-vectors
 175 from the p -th duration group. Eq. 3 becomes DI-PLDA if the SNR-related
 176 term is replaced by a duration-related term, i.e.,

$$\mathbf{x}_{ij}^p = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{R}\mathbf{y}_p + \boldsymbol{\epsilon}_{ij}^p, \quad (4)$$

177 where \mathbf{R} defines the duration subspace, \mathbf{y}_p is a latent duration factor with
 178 a standard normal distribution. Other terms have the same meaning as in
 179 Eq. 3.

180 In [21], the latent factors \mathbf{h}_i and \mathbf{y}_p are assumed to be posteriorly in-
 181 dependent. In this paper, we consider \mathbf{h}_i and \mathbf{y}_p are posteriorly dependent
 182 and use variational Bayes methods [2] to derive EM algorithms for training
 183 the SI-PLDA and DI-PLDA models.

184 Denote $N_i = \sum_{p=1}^P H_i(p)$ as the number of training utterances from the
 185 i -th speaker and $B_p = \sum_{i=1}^S H_i(p)$ as the number of the training utterances
 186 in the p -th duration group. Given an old estimate of the model parameters
 187 $\boldsymbol{\theta} = \{\mathbf{m}, \mathbf{V}, \mathbf{R}, \boldsymbol{\Sigma}\}$, we aim to find a new estimate $\boldsymbol{\theta}'$ that maximizes the
 188 auxiliary function:

$$\begin{aligned} Q(\boldsymbol{\theta}'|\boldsymbol{\theta}) &= \mathbb{E}_{q(\underline{\mathbf{h}}, \underline{\mathbf{y}})} \left\{ \ln p(\mathcal{X}, \underline{\mathbf{h}}, \underline{\mathbf{y}}|\boldsymbol{\theta}') \middle| \mathcal{X}, \boldsymbol{\theta} \right\} \\ &= \mathbb{E}_{q(\underline{\mathbf{h}}, \underline{\mathbf{y}})} \left\{ \sum_{ijp} \ln [p(\mathbf{x}_{ij}^p | \mathbf{h}_i, \mathbf{y}_p, \boldsymbol{\theta}') p(\mathbf{h}_i, \mathbf{y}_p)] \middle| \mathcal{X}, \boldsymbol{\theta} \right\}, \end{aligned} \quad (5)$$

189 where $\underline{\mathbf{h}} = \{\mathbf{h}_1, \dots, \mathbf{h}_s\}$, $\underline{\mathbf{y}} = \{\mathbf{y}_1, \dots, \mathbf{y}_p\}$, and $q(\underline{\mathbf{h}}, \underline{\mathbf{y}})$ is the variational
 190 posterior density of $\underline{\mathbf{h}}$ and $\underline{\mathbf{y}}$. To maximize $Q(\boldsymbol{\theta}'|\boldsymbol{\theta})$, we differentiate $Q(\boldsymbol{\theta}'|\boldsymbol{\theta})$
 191 with respect to the model parameters $\{\mathbf{m}, \mathbf{V}, \mathbf{R}, \boldsymbol{\Sigma}\}$ and set the resulting
 192 derivatives to $\mathbf{0}$. This leads to

$$\mathbf{m} = \frac{1}{N} \sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \mathbf{x}_{ij}^p \quad (6)$$

$$\mathbf{V}' = \left\{ \sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \left[(\mathbf{x}_{ij}^p - \mathbf{m}) \langle \mathbf{h}_i | \mathcal{X} \rangle - \mathbf{R} \langle \mathbf{y}_p \mathbf{h}_i^\top | \mathcal{X} \rangle \right] \right\} \left[\sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \langle \mathbf{h}_i \mathbf{h}_i^\top | \mathcal{X} \rangle \right]^{-1} \quad (7)$$

$$\mathbf{R}' = \left\{ \sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \left[(\mathbf{x}_{ij}^p - \mathbf{m}) \langle \mathbf{y}_p | \mathcal{X} \rangle - \mathbf{V} \langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle \right] \right\} \left[\sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \langle \mathbf{y}_p \mathbf{y}_p^\top | \mathcal{X} \rangle \right]^{-1} \quad (8)$$

$$\mathbf{\Sigma}' = \frac{1}{N} \sum_{i=1}^S \sum_{p=1}^P \sum_{j=1}^{H_i(p)} \left[(\mathbf{x}_{ij}^p - \mathbf{m})(\mathbf{x}_{ij}^p - \mathbf{m})^\top - \mathbf{V} \langle \mathbf{h}_i | \mathcal{X} \rangle (\mathbf{x}_{ij}^p - \mathbf{m})^\top - \mathbf{R} \langle \mathbf{y}_p | \mathcal{X} \rangle (\mathbf{x}_{ij}^p - \mathbf{m})^\top \right] \quad (9)$$

195 where $N = \sum_{i=1}^S N_i = \sum_{p=1}^P B_p$.

196 Eq. 6–Eq. 9 constitute the M-step of the EM algorithm. To update
197 the model parameters in the M-step, we need to estimate the posterior
198 distribution of \mathbf{h}_i and \mathbf{y}_p . These posteriors can be obtained through the
199 variational Bayes method as explained below.

We approximate the true posterior $p(\underline{\mathbf{h}}, \underline{\mathbf{y}} | \mathcal{X})$ by a variational posterior $q(\underline{\mathbf{h}}, \underline{\mathbf{y}})$ and write the marginal likelihood of \mathcal{X} as

$$\begin{aligned} \ln p(\mathcal{X}) &= \int \int q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \ln p(\mathcal{X}) d\underline{\mathbf{h}} d\underline{\mathbf{y}} \\ &= \int \int q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \ln \left[\frac{p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X})}{p(\underline{\mathbf{h}}, \underline{\mathbf{y}} | \mathcal{X})} \right] d\underline{\mathbf{h}} d\underline{\mathbf{y}} \\ &= \int \int q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \ln \left[\frac{p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X})}{q(\underline{\mathbf{h}}, \underline{\mathbf{y}})} \right] d\underline{\mathbf{h}} d\underline{\mathbf{y}} + \int \int q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \ln \left[\frac{q(\underline{\mathbf{h}}, \underline{\mathbf{y}})}{p(\underline{\mathbf{h}}, \underline{\mathbf{y}} | \mathcal{X})} \right] d\underline{\mathbf{h}} d\underline{\mathbf{y}} \\ &= \mathcal{L}(q) + \mathcal{D}_{\text{KL}}(q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \| p(\underline{\mathbf{h}}, \underline{\mathbf{y}} | \mathcal{X})). \end{aligned} \quad (10)$$

200 In Eq. 10, $\mathcal{D}_{\text{KL}}(q \| p)$ is the KL-divergence between distributions q and p and

$$\mathcal{L}(q) = \int \int q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \ln \left[\frac{p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X})}{q(\underline{\mathbf{h}}, \underline{\mathbf{y}})} \right] d\underline{\mathbf{h}} d\underline{\mathbf{y}} \quad (11)$$

201 is the variational lower bound of the marginal likelihood. Since KL-divergence
202 is non-negative, we can maximize the marginal likelihood through maximiz-
203 ing the lower bound with respect to $q(\underline{\mathbf{h}}, \underline{\mathbf{y}})$. The maximum occurs when
204 $q(\underline{\mathbf{h}}, \underline{\mathbf{y}})$ equals the true posterior $p(\underline{\mathbf{h}}, \underline{\mathbf{y}} | \mathcal{X})$. Then, we assume that the ap-

205 proximated posterior $q(\underline{\mathbf{h}}, \underline{\mathbf{y}})$ can be factorized as follows:

$$\ln q(\underline{\mathbf{h}}, \underline{\mathbf{y}}) = \ln q(\underline{\mathbf{h}}) + \ln q(\underline{\mathbf{y}}) = \sum_{i=1}^S \ln q(\mathbf{h}_i) + \sum_{p=1}^P \ln q(\mathbf{y}_p). \quad (12)$$

206 By maximizing the lower bound $\mathcal{L}(q)$ in Eq. 11, we obtain [2, 35]

$$\begin{aligned} \ln q(\underline{\mathbf{h}}) &= \mathbb{E}_{q(\underline{\mathbf{y}})} \{\ln p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X})\} + \text{const} \\ \ln q(\underline{\mathbf{y}}) &= \mathbb{E}_{q(\underline{\mathbf{h}})} \{\ln p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X})\} + \text{const}, \end{aligned} \quad (13)$$

207 where $\mathbb{E}_{q(\underline{\mathbf{y}})}$ means taking expectation with respect to $\underline{\mathbf{y}}$ using $q(\underline{\mathbf{y}})$ as the
208 density.

209 Note that $\ln q(\underline{\mathbf{h}})$ in Eq. 13 can be written as

$$\begin{aligned} \ln q(\underline{\mathbf{h}}) &= \sum_i \ln q(\mathbf{h}_i) = \langle \ln p(\underline{\mathbf{h}}, \underline{\mathbf{y}}, \mathcal{X}) \rangle_{\underline{\mathbf{y}}} + \text{const} \\ &= \langle \ln p(\mathcal{X} | \underline{\mathbf{h}}, \underline{\mathbf{y}}) \rangle_{\underline{\mathbf{y}}} + \langle \ln p(\underline{\mathbf{h}}, \underline{\mathbf{y}}) \rangle_{\underline{\mathbf{y}}} + \text{const} \\ &= \sum_{ijp} \left\langle \ln \mathcal{N}(\mathbf{x}_{ij}^p | \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{R}\mathbf{y}_p, \boldsymbol{\Sigma}) \right\rangle_{\mathbf{y}_p} + \sum_i \langle \ln \mathcal{N}(\mathbf{h}_i | \mathbf{0}, \mathbf{I}) \rangle_{\underline{\mathbf{y}}} \\ &\quad + \sum_p \langle \ln \mathcal{N}(\mathbf{y}_p | \mathbf{0}, \mathbf{I}) \rangle_{\mathbf{y}_p} + \text{const} \\ &= -\frac{1}{2} \sum_{ijp} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{V}\mathbf{h}_i - \mathbf{R}\mathbf{y}_p^*)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{V}\mathbf{h}_i - \mathbf{R}\mathbf{y}_p^*) - \frac{1}{2} \sum_i \mathbf{h}_i^\top \mathbf{h}_i + \text{const}^1 \\ &= \sum_i \left[\mathbf{h}_i^\top \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \sum_{jp} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{R}\mathbf{y}_p^*) - \frac{1}{2} \mathbf{h}_i^\top \left(\mathbf{I} + \sum_p H_i(p) \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V} \right) \mathbf{h}_i \right] + \text{const}. \end{aligned} \quad (14)$$

210 where $\mathbf{y}_p^* \equiv \langle \mathbf{y}_p | \mathcal{X} \rangle_{\mathbf{y}_p}$ is the posterior mean of \mathbf{y}_p in the previous iteration
211 and $\langle \cdot \rangle_{\mathbf{y}_p}$ denotes the expectation with respect to \mathbf{y}_p .

212 By reading off \mathbf{h}_i in Eq. 14 and comparing with $\sum_i \ln q(\mathbf{h}_i)$, we note that
213 $q(\mathbf{h}_i)$ is a Gaussian with the following mean vector and precision matrix:

$$\mathbb{E}_{q(\mathbf{h}_i)} \{\mathbf{h}_i | \mathcal{X}\} = \langle \mathbf{h}_i | \mathcal{X} \rangle = \left(\mathbf{L}_i^{(1)} \right)^{-1} \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \sum_{p=1}^P \sum_{j=1}^{H_i(p)} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{R}\mathbf{y}_p^*)$$

¹ $\langle \ln \mathcal{N}(\mathbf{y}_p | \mathbf{0}, \mathbf{I}) \rangle_{\mathbf{y}_p}$ is the differential entropy of normal distribution and is independent of \mathbf{h}_i , see Chapter 8 in [28].

214 and

$$\mathbf{L}_i^{(1)} \equiv \mathbf{I} + \sum_{p=1}^P H_i(p) \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V}.$$

215 As a result, the second-order moment required in the M-step can be com-
216 puted as follows:

$$\langle \mathbf{h}_i \mathbf{h}_i^\top | \mathcal{X} \rangle = \left(\mathbf{L}_i^{(1)} \right)^{-1} + \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top.$$

217 Similarly, the posterior mean and second-order moment of \mathbf{y}_p can also be
218 obtained by comparing the terms in $\ln q(\mathbf{y}_p)$ with a Gaussian distribution.

219 The M-step also requires the posterior moment $\langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle$, which can be
220 approximated by using variational Bayes principle:

$$p(\mathbf{h}_i, \mathbf{y}_p | \mathcal{X}) \approx q(\mathbf{h}_i) q(\mathbf{y}_p), \quad (15)$$

221 where both $q(\mathbf{h}_i)$ and $q(\mathbf{y}_p)$ are Gaussians. Based on the law of total expect-
222 ation [1], the factorization in Eq. 15 gives

$$\begin{aligned} \langle \mathbf{y}_p \mathbf{h}_i^\top | \mathcal{X} \rangle &\approx \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \\ \langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle &\approx \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top. \end{aligned}$$

223 Therefore, the equations for the variational E-step are as follows:

$$\mathbf{L}_i^{(1)} = \mathbf{I} + N_i \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V} \quad i = 1, \dots, S \quad (16)$$

224

$$\mathbf{L}_p^{(2)} = \mathbf{I} + B_p \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \mathbf{R} \quad p = 1, \dots, P \quad (17)$$

225

$$\langle \mathbf{h}_i | \mathcal{X} \rangle = \left(\mathbf{L}_i^{(1)} \right)^{-1} \mathbf{V}^\top \boldsymbol{\Sigma}^{-1} \sum_{p=1}^P \sum_{j=1}^{H_i(p)} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{R} \mathbf{y}_p^*) \quad (18)$$

226

$$\langle \mathbf{y}_p | \mathcal{X} \rangle = \left(\mathbf{L}_p^{(2)} \right)^{-1} \mathbf{R}^\top \boldsymbol{\Sigma}^{-1} \sum_{i=1}^S \sum_{j=1}^{H_i(p)} (\mathbf{x}_{ij}^p - \mathbf{m} - \mathbf{V} \mathbf{h}_i^*) \quad (19)$$

227

$$\langle \mathbf{h}_i \mathbf{h}_i^\top | \mathcal{X} \rangle = \left(\mathbf{L}_i^{(1)} \right)^{-1} + \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \quad (20)$$

228

$$\langle \mathbf{y}_p \mathbf{y}_p^\top | \mathcal{X} \rangle = \left(\mathbf{L}_p^{(2)} \right)^{-1} + \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top \quad (21)$$

229

$$\langle \mathbf{y}_p \mathbf{h}_i^\top | \mathcal{X} \rangle \approx \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \quad (22)$$

230

$$\langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle \approx \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top \quad (23)$$

231 Algorithm 1 shows the procedure of training a duration-invariant PLDA
 232 model.

Algorithm 1 Variational Bayes EM Algorithm for Duration-Invariant PLDA

Input:

Development data set consisting of i-vectors $\mathcal{X} = \{\mathbf{x}_{ij}^p | i = 1, \dots, S; j = 1, \dots, H_i(p); p = 1, \dots, P\}$, with identity labels and duration group labels.

Initialization:

$\mathbf{y}_p^* \leftarrow \mathbf{0}$;

$\Sigma \leftarrow 0.01\mathbf{I}$;

$\mathbf{V}, \mathbf{R} \leftarrow$ eigenvectors of PCA projection matrix learned using data set \mathcal{X} ;

Parameter Estimation:

- 1) Compute \mathbf{m} via Eq. 6;
- 2) Compute $\mathbf{L}_i^{(1)}$ and $\mathbf{L}_p^{(2)}$ according to Eq. 16 and Eq. 17, respectively;
- 3) Set \mathbf{y}_p^* to the posterior mean of \mathbf{y}_p . Compute the posterior mean of \mathbf{h}_i using Eq. 18;
- 4) Use the posterior mean of \mathbf{h}_i computed in Step 3 to update the posterior mean of \mathbf{y}_p according to Eq. 19;
- 5) Compute the other terms in the E-step (Eq. 20–Eq. 23);
- 6) Update the model parameters using Eq. 7 to Eq. 9;
- 7) Go to Step 2 until convergence;

Return: the parameters of the duration-invariant PLDA model $\theta = \{\mathbf{m}, \mathbf{V}, \mathbf{R}, \Sigma\}$.

232

233 *3.2. Likelihood Ratio Scores*

234 If the durations of target and test utterances are not known (or not
 235 used), the likelihood ratio score can be computed in the same manner as
 236 in SI-PLDA [21]. Because the duration ℓ is usually known in practice, the
 237 likelihood ratio score can be also computed as follows:

$$S_{\text{LR}}(\mathbf{x}_s, \mathbf{x}_t | \ell_s, \ell_t) = \ln \frac{p(\mathbf{x}_s, \mathbf{x}_t | \text{same-speaker}, \ell_s, \ell_t)}{p(\mathbf{x}_s, \mathbf{x}_t | \text{different-speakers}, \ell_s, \ell_t)}, \quad (24)$$

238 where \mathbf{x}_s and \mathbf{x}_t denote the target-speaker’s i-vector and test i-vector, respec-
 239 tively, and ℓ_s and ℓ_t denote the durations of the corresponding utterances.

240 Based on different assumptions on the posterior density of \mathbf{y}_p , we pro-
 241 pose two methods to calculate the score. They are derived in the following
 242 subsections.

243 *3.2.1. Duration Factors with a Sharp Posterior Density*

244 Assume that the duration ℓ of an utterance belongs to the p -th duration
 245 group and that the posterior density of \mathbf{y}_p is sharp at its mean \mathbf{y}_p^* .² Then,
 246 the marginal-likelihood of i-vector \mathbf{x} can be written as:

$$\begin{aligned}
 p(\mathbf{x}|\ell \in p\text{-th duration group}) &= \int_{\mathbf{h}} p(\mathbf{x}|\mathbf{h}, \mathbf{y}_p^*) p(\mathbf{h}) d\mathbf{h} \\
 &= \int_{\mathbf{h}} \mathcal{N}(\mathbf{x}|\mathbf{m} + \mathbf{V}\mathbf{h} + \mathbf{R}\mathbf{y}_p^*, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{h}|\mathbf{0}, \mathbf{I}) d\mathbf{h} \\
 &= \mathcal{N}(\mathbf{x}|\mathbf{m} + \mathbf{R}\mathbf{y}_p^*, \mathbf{V}\mathbf{V}^\top + \boldsymbol{\Sigma}),
 \end{aligned} \tag{25}$$

247 where $\mathbf{y}_p^* \equiv \langle \mathbf{y}_p | \mathcal{X} \rangle$ is the posterior mean of \mathbf{y}_p . Given a test i-vector \mathbf{x}_t and
 248 a target i-vector \mathbf{x}_s , we can use Eq. 25 to compute the likelihood ratio score:

$$\begin{aligned}
 S_{\text{LR}}(\mathbf{x}_s, \mathbf{x}_t | \ell_s, \ell_t) &= \ln \frac{p(\mathbf{x}_s, \mathbf{x}_t | \text{same-speaker}, \ell_s, \ell_t)}{p(\mathbf{x}_s, \mathbf{x}_t | \text{different-speakers}, \ell_s, \ell_t)} \\
 &= \ln \frac{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} + \mathbf{R}\mathbf{y}_{p_s}^* \\ \mathbf{m} + \mathbf{R}\mathbf{y}_{p_t}^* \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Psi} & \boldsymbol{\Sigma}_{ac} \\ \boldsymbol{\Sigma}_{ac} & \boldsymbol{\Psi} \end{bmatrix}\right)}{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} + \mathbf{R}\mathbf{y}_{p_s}^* \\ \mathbf{m} + \mathbf{R}\mathbf{y}_{p_t}^* \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Psi} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Psi} \end{bmatrix}\right)} \\
 &= \frac{1}{2} [\bar{\mathbf{x}}_s^\top \mathbf{Q} \bar{\mathbf{x}}_s + 2\bar{\mathbf{x}}_s^\top \mathbf{P} \bar{\mathbf{x}}_t + \bar{\mathbf{x}}_t^\top \mathbf{Q} \bar{\mathbf{x}}_t] + \text{const}
 \end{aligned} \tag{26}$$

where

$$\begin{aligned}
 \bar{\mathbf{x}}_s &= \mathbf{x}_s - \mathbf{m} - \mathbf{R}\mathbf{y}_{p_s}^* \\
 \bar{\mathbf{x}}_t &= \mathbf{x}_t - \mathbf{m} - \mathbf{R}\mathbf{y}_{p_t}^* \\
 \mathbf{Q} &= \boldsymbol{\Psi}^{-1} - (\boldsymbol{\Psi} - \boldsymbol{\Sigma}_{ac} \boldsymbol{\Psi}^{-1} \boldsymbol{\Sigma}_{ac})^{-1} \\
 \mathbf{P} &= \boldsymbol{\Psi}^{-1} \boldsymbol{\Sigma}_{ac} (\boldsymbol{\Psi} - \boldsymbol{\Sigma}_{ac} \boldsymbol{\Psi}^{-1} \boldsymbol{\Sigma}_{ac})^{-1} \\
 \boldsymbol{\Psi} &= \mathbf{V}\mathbf{V}^\top + \boldsymbol{\Sigma}; \quad \boldsymbol{\Sigma}_{ac} = \mathbf{V}\mathbf{V}^\top.
 \end{aligned}$$

249 *3.2.2. Duration Factors with a Moderately Sharp Posterior*

250 If the duration ℓ of an utterance falls on the p -th duration group and
 251 the posterior density of \mathbf{y}_p is moderately sharp and follows a Gaussian

²This occurs when the number of training i-vectors B_p in the p -th duration group is large, as suggested by Eq. 17.

252 $\mathcal{N}(\mathbf{y}_p | \boldsymbol{\mu}_p^*, \boldsymbol{\Sigma}_p^*)$,³ the marginal-likelihood of i-vector \mathbf{x} is:

$$\begin{aligned}
p(\mathbf{x} | \ell \in p\text{-th duration group}) &= \int_{\mathbf{h}} \int_{\mathbf{y}_p} p(\mathbf{x} | \mathbf{h}, \mathbf{y}_p) p(\mathbf{h}) p(\mathbf{y}_p) d\mathbf{h} d\mathbf{y}_p \\
&= \int_{\mathbf{h}} \int_{\mathbf{y}_p} \mathcal{N}(\mathbf{x} | \mathbf{m} + \mathbf{V}\mathbf{h} + \mathbf{R}\mathbf{y}_p, \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{h} | \mathbf{0}, \mathbf{I}) \mathcal{N}(\mathbf{y}_p | \mathbf{0}, \mathbf{I}) d\mathbf{h} d\mathbf{y}_p \\
&= \int_{\mathbf{y}_p} \mathcal{N}(\mathbf{x} | \mathbf{m} + \mathbf{R}\mathbf{y}_p, \mathbf{V}\mathbf{V}^\top + \boldsymbol{\Sigma}) \mathcal{N}(\mathbf{y}_p | \mathbf{0}, \mathbf{I}) d\mathbf{y}_p \\
&= \mathcal{N}(\mathbf{x} | \mathbf{m} + \mathbf{R}\boldsymbol{\mu}_p^*, \mathbf{V}\mathbf{V}^\top + \mathbf{R}\boldsymbol{\Sigma}_p^*\mathbf{R}^\top + \boldsymbol{\Sigma}),
\end{aligned} \tag{27}$$

253 where $\boldsymbol{\mu}_p^*$ can be computed according to Eq. 19 and $\boldsymbol{\Sigma}_p^*$ can be estimated
254 from the inverse of $\mathbf{L}_p^{(2)}$ in Eq. 17. Given a test i-vector \mathbf{x}_t and a target
255 i-vector \mathbf{x}_s , the likelihood ratio score can be computed as:

$$\begin{aligned}
S_{\text{LR}}(\mathbf{x}_s, \mathbf{x}_t | \ell_s, \ell_t) &= \ln \frac{p(\mathbf{x}_s, \mathbf{x}_t | \text{same-speaker}, \ell_s, \ell_t)}{p(\mathbf{x}_s, \mathbf{x}_t | \text{different-speakers}, \ell_s, \ell_t)} \\
&= \ln \frac{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} + \mathbf{R}\boldsymbol{\mu}_{p_s}^* \\ \mathbf{m} + \mathbf{R}\boldsymbol{\mu}_{p_t}^* \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_s & \boldsymbol{\Sigma}_{ac} \\ \boldsymbol{\Sigma}_{ac} & \boldsymbol{\Sigma}_t \end{bmatrix}\right)}{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} + \mathbf{R}\boldsymbol{\mu}_{p_s}^* \\ \mathbf{m} + \mathbf{R}\boldsymbol{\mu}_{p_t}^* \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_s & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_t \end{bmatrix}\right)} \\
&= \frac{1}{2} [\bar{\mathbf{x}}_s^\top \mathbf{A}_{s,t} \bar{\mathbf{x}}_s + 2\bar{\mathbf{x}}_s^\top \mathbf{B}_{s,t} \bar{\mathbf{x}}_t + \bar{\mathbf{x}}_t^\top \mathbf{C}_{s,t} \bar{\mathbf{x}}_t] + \text{const}
\end{aligned} \tag{28}$$

where

$$\begin{aligned}
\bar{\mathbf{x}}_s &= \mathbf{x}_s - \mathbf{m} - \mathbf{R}\boldsymbol{\mu}_{p_s}^*; & \bar{\mathbf{x}}_t &= \mathbf{x}_t - \mathbf{m} - \mathbf{R}\boldsymbol{\mu}_{p_t}^* \\
\mathbf{A}_{s,t} &= \boldsymbol{\Sigma}_s^{-1} - (\boldsymbol{\Sigma}_s - \boldsymbol{\Sigma}_{ac}\boldsymbol{\Sigma}_t^{-1}\boldsymbol{\Sigma}_{ac})^{-1} \\
\mathbf{B}_{s,t} &= \boldsymbol{\Sigma}_s^{-1}\boldsymbol{\Sigma}_{ac}(\boldsymbol{\Sigma}_t - \boldsymbol{\Sigma}_{ac}\boldsymbol{\Sigma}_s^{-1}\boldsymbol{\Sigma}_{ac})^{-1} \\
\mathbf{C}_{s,t} &= \boldsymbol{\Sigma}_t^{-1} - (\boldsymbol{\Sigma}_t - \boldsymbol{\Sigma}_{ac}\boldsymbol{\Sigma}_s^{-1}\boldsymbol{\Sigma}_{ac})^{-1} \\
\boldsymbol{\Sigma}_s &= \mathbf{V}\mathbf{V}^\top + \mathbf{R}\boldsymbol{\Sigma}_{p_s}^*\mathbf{R}^\top + \boldsymbol{\Sigma}; & \boldsymbol{\Sigma}_t &= \mathbf{V}\mathbf{V}^\top + \mathbf{R}\boldsymbol{\Sigma}_{p_t}^*\mathbf{R}^\top + \boldsymbol{\Sigma}; & \boldsymbol{\Sigma}_{ac} &= \mathbf{V}\mathbf{V}^\top.
\end{aligned}$$

256 4. SNR- and Duration-invariant PLDA

257 This section describes a new modeling approach, namely SNR- and
258 duration-invariant PLDA (SDI-PLDA), for robust speaker verification. Un-

³This occurs when the number of training i-vectors B_p in the p -th duration group is moderate, as suggested by Eq. 17.

259 like conventional Gaussian PLDA and SNR-invariant PLDA, the proposed
 260 model has three labeled latent factors representing speaker-specific, SNR-
 261 specific and duration-specific information, respectively.

262 4.1. Generative Model

263 The SNR- and duration-invariant PLDA is inspired by the notion of
 264 Gaussian PLDA in which i-vectors from the same speaker should share a
 265 speaker latent factor. Similarly, this method is based on two hypotheses:
 266 (1) i-vectors derived from utterances that fall within a narrow SNR range
 267 should share similar SNR-specific information; and (2) i-vectors extracted
 268 from utterances with comparable durations should share similar duration-
 269 specific information.

270 To confirm the first hypothesis, we plotted three groups of i-vectors on
 271 the first 3 principal components in Fig. 2(a), where each group corresponds
 272 to a specific SNR-level shown in the legend. To ensure that the cluster
 273 displacement is not caused by speaker variability, each group contains the
 274 i-vectors from the same set of speakers. Evidently, the i-vectors form three
 275 clusters, one for each SNR group. To illustrate the the second hypothesis,
 276 we display three groups of i-vectors on their first 3 principal components
 277 in Fig. 2(b), where each group corresponds to one duration range shown in
 278 the legend. To ensure that the variability in i-vectors is not due to noise-
 279 level and speaker variabilities, all of the i-vectors were obtained from clean
 280 telephone conversations and each duration group comprises the same set of
 281 target speakers. Evidently, the i-vectors form three clusters and the locations
 282 of the clusters depend on the duration range.

283 From a modeling standpoint, both SNR-specific and duration-specific
 284 information can be captured using latent factors just like speaker factor in
 285 conventional PLDA model. We refer to these latent factors as SNR factor
 286 and duration factor in the remainder of this paper.

287 Under the above assumptions, an LDA- or NFA-projected i-vector [21]
 288 can be regarded as an observation generated from a linear generative model
 289 that comprises four components: (1) speaker component, (2) SNR compo-
 290 nent, (3) duration component, and (4) channel variability and the remaining
 291 variability that cannot be captured by the first three components. Assume
 292 that we have a set of D -dimensional NFA-projected i-vectors $\mathcal{X} = \{\hat{\mathbf{x}}_{ij}^{kp} | i =$
 293 $1, \dots, S; k = 1, \dots, K; p = 1, \dots, P; j = 1, \dots, H_{ik}(p)\}$ obtained from S
 294 speakers, where $\hat{\mathbf{x}}_{ij}^{kp}$ is the j -th i-vector from speaker i , and k and p index to
 295 the SNR and duration groups to which the corresponding utterances belong,

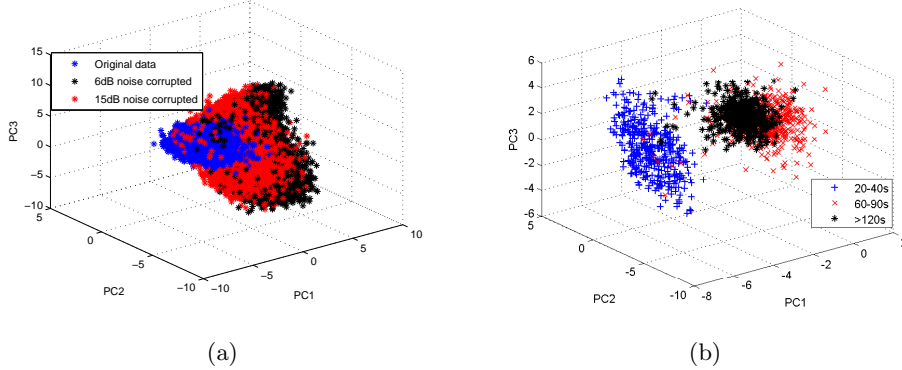


Figure 2: (a) Projection of i-vectors derived from utterances with different SNR-levels on their first three principal components. (b) Projection of i-vectors derived from variable-length utterances on their first three principal components.

296 respectively. In the proposed model, $\hat{\mathbf{x}}_{ij}^{kp}$ can be expressed as:

$$\hat{\mathbf{x}}_{ij}^{kp} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{U}\mathbf{w}_k + \mathbf{R}\mathbf{y}_p + \boldsymbol{\epsilon}_{ij}^{kp}, \quad (29)$$

297 where \mathbf{m} is a $D \times 1$ vector representing the global offset, \mathbf{h}_i is a $Q_1 \times 1$ vector
 298 denoting the speaker factor with prior distribution $\mathcal{N}(\mathbf{h}|\mathbf{0}, \mathbf{I})$, \mathbf{w}_k is a $Q_2 \times 1$
 299 vector denoting the latent SNR factor with a prior distribution of $\mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{I})$,
 300 \mathbf{y}_p is a $Q_3 \times 1$ vector denoting the latent duration factor with a standard
 301 normal prior, $\boldsymbol{\epsilon}_{ij}^{kp}$ is a $D \times 1$ vector denoting the residue which follows a
 302 Gaussian distribution $\mathcal{N}(\boldsymbol{\epsilon}|\mathbf{0}, \boldsymbol{\Sigma})$, \mathbf{V} is a $D \times Q_1$ matrix whose columns span
 303 the speaker subspace, \mathbf{U} is a $D \times Q_2$ matrix whose columns span the SNR
 304 subspace, and \mathbf{R} is a $D \times Q_3$ matrix which defines the duration subspace.
 305 \mathbf{h}_i , \mathbf{w}_k , and \mathbf{y}_p are assumed to be independent in their prior. Fig. 3 shows
 306 the graphical model of SDI-PLDA.

307 The proposed SNR- and duration-invariant PLDA is different from the
 308 conventional PLDA in that the former makes use of multiple labels (speaker
 309 IDs, SNR levels, and duration ranges) for training the loading matrices,
 310 whereas the latter only uses the speaker IDs. To exploit the duration in-
 311 formation in the training utterances, the proposed model has an additional
 312 subspace called duration subspace, which results in an extra latent factor
 313 called duration factor. Unlike the term $\mathbf{G}\mathbf{r}_{ij}$ in Eq. 1, which is speaker- and
 314 session-dependent, the SNR component $\mathbf{U}\mathbf{w}_k$ and the duration component
 315 $\mathbf{R}\mathbf{y}_p$ in Eq. 29 depend on the SNR groups and duration groups, respectively.

331

$$\langle \mathbf{h}_i | \mathcal{X} \rangle = (\mathbf{L}_i^{(1)})^{-1} \mathbf{V}^\top \Sigma^{-1} \sum_{k p j} (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m} - \mathbf{U} \mathbf{w}_k^* - \mathbf{R} \mathbf{y}_p^*) \quad (34)$$

332

$$\langle \mathbf{w}_k | \mathcal{X} \rangle = (\mathbf{L}_k^{(2)})^{-1} \mathbf{U}^\top \Sigma^{-1} \sum_{i p j} (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m} - \mathbf{V} \mathbf{h}_i^* - \mathbf{R} \mathbf{y}_p^*) \quad (35)$$

333

$$\langle \mathbf{y}_p | \mathcal{X} \rangle = (\mathbf{L}_p^{(3)})^{-1} \mathbf{R}^\top \Sigma^{-1} \sum_{i k j} (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m} - \mathbf{V} \mathbf{h}_i^* - \mathbf{U} \mathbf{w}_k^*) \quad (36)$$

334

$$\langle \mathbf{h}_i \mathbf{h}_i^\top | \mathcal{X} \rangle = (\mathbf{L}_i^{(1)})^{-1} + \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \quad (37)$$

335

$$\langle \mathbf{w}_k \mathbf{w}_k^\top | \mathcal{X} \rangle = (\mathbf{L}_k^{(2)})^{-1} + \langle \mathbf{w}_k | \mathcal{X} \rangle \langle \mathbf{w}_k | \mathcal{X} \rangle^\top \quad (38)$$

336

$$\langle \mathbf{y}_p \mathbf{y}_p^\top | \mathcal{X} \rangle = (\mathbf{L}_p^{(3)})^{-1} + \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top \quad (39)$$

337

$$\langle \mathbf{w}_k \mathbf{h}_i^\top | \mathcal{X} \rangle \approx \langle \mathbf{w}_k | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \quad (40)$$

338

$$\langle \mathbf{h}_i \mathbf{w}_k^\top | \mathcal{X} \rangle \approx \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{w}_k | \mathcal{X} \rangle^\top \quad (41)$$

339

$$\langle \mathbf{w}_k \mathbf{y}_p^\top | \mathcal{X} \rangle \approx \langle \mathbf{w}_k | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top \quad (42)$$

340

$$\langle \mathbf{y}_p \mathbf{w}_k^\top | \mathcal{X} \rangle \approx \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{w}_k | \mathcal{X} \rangle^\top \quad (43)$$

341

$$\langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle \approx \langle \mathbf{h}_i | \mathcal{X} \rangle \langle \mathbf{y}_p | \mathcal{X} \rangle^\top \quad (44)$$

342

$$\langle \mathbf{y}_p \mathbf{h}_i^\top | \mathcal{X} \rangle \approx \langle \mathbf{y}_p | \mathcal{X} \rangle \langle \mathbf{h}_i | \mathcal{X} \rangle^\top \quad (45)$$

343 where \mathbf{w}_k^* , \mathbf{y}_p^* , and \mathbf{h}_i^* denote the posterior mean of \mathbf{w}_k , \mathbf{y}_p , and \mathbf{h}_i in the
344 previous iteration, respectively.

345 Given Eq. 31–Eq. 45, the model parameters $\boldsymbol{\theta}'$ can be estimated via the
346 M-step is as follows:

$$\mathbf{m} = \frac{1}{N} \sum_{i k p j} \hat{\mathbf{x}}_{ij}^{kp} \quad (46)$$

347

$$\mathbf{V}' = \left\{ \sum_{i k p j} \left[(\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m}) \langle \mathbf{h}_i | \mathcal{X} \rangle^\top - \mathbf{U} \langle \mathbf{w}_k \mathbf{h}_i^\top | \mathcal{X} \rangle - \mathbf{R} \langle \mathbf{y}_p \mathbf{h}_i^\top | \mathcal{X} \rangle \right] \right\} \left[\sum_{i k p j} \langle \mathbf{h}_i \mathbf{h}_i^\top | \mathcal{X} \rangle \right]^{-1} \quad (47)$$

348

$$\mathbf{U}' = \left\{ \sum_{i k p j} \left[(\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m}) \langle \mathbf{w}_k | \mathcal{X} \rangle^\top - \mathbf{V} \langle \mathbf{h}_i \mathbf{w}_k^\top | \mathcal{X} \rangle - \mathbf{R} \langle \mathbf{y}_p \mathbf{w}_k^\top | \mathcal{X} \rangle \right] \right\} \left[\sum_{i k p j} \langle \mathbf{w}_k \mathbf{w}_k^\top | \mathcal{X} \rangle \right]^{-1} \quad (48)$$

349

$$\mathbf{R}' = \left\{ \sum_{i k p j} \left[(\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m}) \langle \mathbf{y}_p | \mathcal{X} \rangle^\top - \mathbf{V} \langle \mathbf{h}_i \mathbf{y}_p^\top | \mathcal{X} \rangle - \mathbf{U} \langle \mathbf{w}_k \mathbf{y}_p^\top | \mathcal{X} \rangle \right] \right\} \left[\sum_{i k p j} \langle \mathbf{y}_p \mathbf{y}_p^\top | \mathcal{X} \rangle \right]^{-1} \quad (49)$$

350

$$\begin{aligned} \Sigma' = \frac{1}{N} \sum_{ikpj} & \left[(\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m})(\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m})^\top - \mathbf{V} \langle \mathbf{h}_i | \mathcal{X} \rangle (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m})^\top \right. \\ & \left. - \mathbf{U} \langle \mathbf{w}_k | \mathcal{X} \rangle (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m})^\top - \mathbf{R} \langle \mathbf{y}_p | \mathcal{X} \rangle (\hat{\mathbf{x}}_{ij}^{kp} - \mathbf{m})^\top \right] \end{aligned} \quad (50)$$

351 where $N = \sum_{i=1}^S N_i = \sum_{k=1}^K M_k$. Algorithm 2 shows the procedure of
 352 applying the variational EM algorithm.

Algorithm 2 Variational Bayes EM Algorithm for SNR- and Duration-Invariant PLDA

Input:

Development data set comprising NFA-projected i-vectors $\mathcal{X} = \{\hat{\mathbf{x}}_{ij}^{kp} | i = 1, \dots, S; k = 1, \dots, K; p = 1, \dots, P; j = 1, \dots, H_{ik}(p)\}$, with speaker labels, SNR group labels, and duration labels.

Initialization:

$\mathbf{y}_p^* \leftarrow \mathbf{0}, \mathbf{w}_k^* \leftarrow \mathbf{0};$

$\Sigma \leftarrow 0.01\mathbf{I};$

$\mathbf{V}, \mathbf{U}, \mathbf{R} \leftarrow$ eigenvectors obtained from the PCA of $\mathcal{X};$

Parameter Estimation:

- 1) Compute \mathbf{m} via Eq. 46;
- 2) Compute $\mathbf{L}_i^{(1)}, \mathbf{L}_k^{(2)}$, and $\mathbf{L}_p^{(3)}$ according to Eq. 31 to Eq. 33, respectively;
- 3) Compute the posterior mean of \mathbf{h}_i using Eq. 34;
- 4) Use the posterior mean of \mathbf{h}_i computed in Step 3 to update the posterior means of \mathbf{w}_k and \mathbf{y}_p using Eq. 35–Eq. 36;
- 5) Compute the other terms in the E-step (Eq. 37–Eq. 45);
- 6) Update the model parameters using Eq. 47 to Eq. 50;
- 7) Set $\mathbf{y}_p^* = \langle \mathbf{y}_p | \mathcal{X} \rangle$, $\mathbf{w}_k^* = \langle \mathbf{w}_k | \mathcal{X} \rangle$, and $\mathbf{h}_i^* = \langle \mathbf{h}_i | \mathcal{X} \rangle$;
- 8) Go to step 2 until convergence;

Return: the parameters of the SNR- and duration-invariant PLDA model $\theta = \{\mathbf{m}, \mathbf{V}, \mathbf{U}, \mathbf{R}, \Sigma\}$.

353 *4.3. Likelihood Ratio Scores*

354 Assume that both the duration and SNR of target-speaker’s utterance
 355 and test utterance are not known. Denote \mathbf{x}_s and \mathbf{x}_t as the NFA-project i-
 356 vectors of the target-speaker and test utterance, respectively, the likelihood

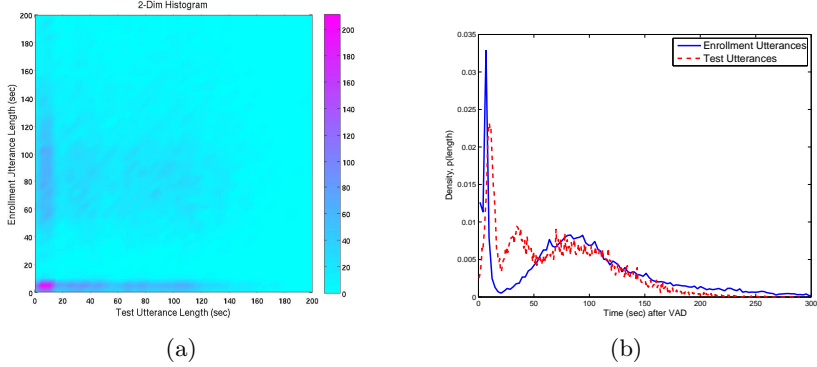


Figure 4: Distribution of utterance duration in NIST 2012 SRE. (a) 2-D histogram showing the length distribution of all possible target-test pairs. (b) Length distributions of enrollment utterances and test utterances (after VAD).

357 ratio score is

$$\begin{aligned}
 S_{\text{LR}}(\mathbf{x}_s, \mathbf{x}_t) &= \ln \frac{P(\mathbf{x}_s, \mathbf{x}_t | \text{same-speaker})}{P(\mathbf{x}_s, \mathbf{x}_t | \text{different-speakers})} \\
 &= \text{const} + \frac{1}{2} \bar{\mathbf{x}}_s^\top \mathbf{Q} \bar{\mathbf{x}}_s + \frac{1}{2} \bar{\mathbf{x}}_t^\top \mathbf{Q} \bar{\mathbf{x}}_t + \bar{\mathbf{x}}_s^\top \mathbf{P} \bar{\mathbf{x}}_t
 \end{aligned} \tag{51}$$

358 where

$$\begin{aligned}
 \bar{\mathbf{x}}_s &= \mathbf{x}_s - \mathbf{m}, & \bar{\mathbf{x}}_t &= \mathbf{x}_t - \mathbf{m}, \\
 \mathbf{P} &= \boldsymbol{\Sigma}_{\text{tot}}^{-1} \boldsymbol{\Sigma}_{ac} (\boldsymbol{\Sigma}_{\text{tot}} - \boldsymbol{\Sigma}_{ac} \boldsymbol{\Sigma}_{\text{tot}}^{-1} \boldsymbol{\Sigma}_{ac})^{-1}, \\
 \mathbf{Q} &= \boldsymbol{\Sigma}_{\text{tot}}^{-1} - (\boldsymbol{\Sigma}_{\text{tot}} - \boldsymbol{\Sigma}_{ac} \boldsymbol{\Sigma}_{\text{tot}}^{-1} \boldsymbol{\Sigma}_{ac})^{-1}, \\
 \boldsymbol{\Sigma}_{ac} &= \mathbf{V}\mathbf{V}^\top, \quad \text{and} \quad \boldsymbol{\Sigma}_{\text{tot}} = \mathbf{V}\mathbf{V}^\top + \mathbf{U}\mathbf{U}^\top + \mathbf{R}\mathbf{R}^\top + \boldsymbol{\Sigma}.
 \end{aligned}$$

359 See Appendix A for the derivation of Eq. 51. When the utterance duration
 360 and SNR are known, the scoring function can be derived using the principles
 361 in Section 3.2.

362 Because \mathbf{P} and \mathbf{Q} can be computed in advance, the computational com-
 363 plexity of SDI-PLDA is the same as that of Gaussian PLDA [7].

364 5. Experimental Setup

365 5.1. Evaluation Protocol and Speech Data

366 Experiments were performed on common conditions (CC) 1 and 4 of the
 367 core set of NIST 2012 Speaker Recognition Evaluation (SRE) [27]. We used

Table 1: Abbreviations of various PLDA models.

Abbreviation	Model Name	Formula
PLDA	Probabilistic LDA	$\mathbf{x}_{ij} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \epsilon_{ij}$ (Eq. 2)
UP-PLDA	Uncertainty propagation PLDA	$\mathbf{x}_{ij} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{R}_{ij}\mathbf{y}_{ij} + \epsilon_{ij}$ (Eq. 2 in [19])
SI-PLDA	SNR-invariant PLDA	$\mathbf{x}_{ij}^k = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{U}\mathbf{w}_k + \epsilon_{ij}^k$ (Eq. 3)
DI-PLDA	Duration-invariant PLDA	$\mathbf{x}_{ij}^p = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{R}\mathbf{y}_p + \epsilon_{ij}^p$ (Eq. 4)
SDI-PLDA	SNR- and duration-invariant PLDA	$\mathbf{x}_{ij}^{kp} = \mathbf{m} + \mathbf{V}\mathbf{h}_i + \mathbf{U}\mathbf{w}_k + \mathbf{R}\mathbf{y}_p + \epsilon_{ij}^{kp}$ (Eq. 29)

368 data from NIST 2005–2010 for system development. The speech data were
 369 divided into the following parts:

- 370 • *Test Data*: Test utterances involved in CC1 comprise clean interview
 371 conversations. Test data in CC4 comprise noise contaminated tele-
 372 phone conversations with SNR ranging from 0dB to 50dB. Readers
 373 may refer to [21] for the SNR distributions of test utterances in these
 374 common conditions and the procedure for measuring SNRs.
- 375 • *Enrollment Data*: Enrollment data for CC1 comprise target-speakers’
 376 conversations recorded by using different types of microphones. For
 377 CC4, enrollment data comprises the telephone conversations of target
 378 speakers. Each target speaker has one or more conversations recorded
 379 over different telephone channels and with different durations longer
 380 than 10 seconds.
- 381 • *Development Data*: Development data were used for estimating the
 382 subspace projection matrices (WCCN and NFA) for i-vector preprocess-
 383 ing. They were also used for estimating the parameters of PLDA,
 384 uncertainty propagation PLDA in [19], SNR-invariant PLDA and SNR-
 385 and duration-invariant PLDA models (see Table 1 for the abbrevia-
 386 tions of these models). For experiments on CC1, the development
 387 data consist of microphone segments in 2005–2010 SREs. For CC4,
 388 the development data comprise two parts. The first part was extracted
 389 from the telephone and microphone segments in 2005–2010 SREs and
 390 the second part was obtained by adding babble noise to the telephone
 391 segments of 2005–2010 SREs at different SNRs. The procedure of pro-
 392 ducing these noisy speech files is described in Section IV-B of [21]. For
 393 each gender, 14,000 noise corrupted files with SNR ranging from 2dB

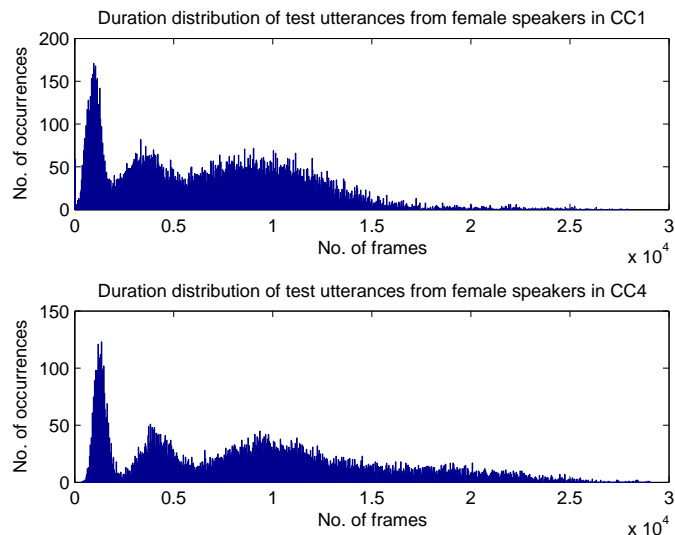


Figure 5: Duration distributions of test utterances (after VAD) from female speakers in CC1 and CC4 of NIST 2012 SRE.

394 to 15dB were randomly selected from all of the noise corrupted files.
 395 Speakers with less than 10 conversations were excluded. Both the
 396 microphone and telephone conversations from NIST 2005–2008 SREs
 397 were used as development data to train the gender-dependent UBMs
 398 and total variability matrices.

399 Fig. 4 shows the duration distribution of the enrollment and test utter-
 400 ances (after VAD) in NIST 2012 SRE. Evidently, there are many trials that
 401 involve short test utterances tested against long enrollment utterances or
 402 long test utterances tested against short enrollment utterances. The dura-
 403 tion distributions of test utterances (after VAD) in CC1 and CC4 are shown
 404 in Fig. 5. It is obvious that the test utterances in CC1 and CC4 covers a
 405 wide range of durations.

406 5.2. Acoustic Feature Extraction

407 For each conversation, a two-channel voice activity detector (VAD) [25,
 408 37] was applied to prune out silence regions. The VAD is specifically de-
 409 signed for NIST SREs. Special attention has been paid to address utterances
 410 with low SNR, impulsive noise, and cross talks in the interview speech files.
 411 The main idea is to apply speech enhancement as a pre-processing step to

412 boost energy contrast between speech and non-speech regions, which facil-
413 itates the subsequence speech/non-speech decisions either by log-likelihood
414 ratio tests or by comparing with energy-based thresholds.

415 The speech regions of each utterance were segmented into 25-ms Ham-
416 ming windowed frames with 10-ms frame shift. For each frame, the first 19
417 Mel frequency cepstral coefficients (MFCC) and log energy together with
418 their first and second derivatives were packed to form a 60-dimensional
419 acoustic vector. Cepstral mean normalization and feature warping [29] with
420 a window size of 3 seconds were then applied to the acoustic vectors.

421 5.3. *I-vector Extraction and PLDA Modeling*

422 I-vectors were extracted based on gender-dependent UBMs with 1024
423 mixtures and total variability matrices with 500 total factors. Similar to [26],
424 we applied within-class covariance normalization (WCCN) [12] to whiten
425 the i-vectors, followed by length normalization (LN) to reduce the non-
426 Gaussian behavior of the 500-dimensional i-vectors. Then, nonparametric
427 feature analysis (NFA) [21, 22] was applied to reduce intra-speaker vari-
428 ability and emphasize discriminative information. This procedure projects
429 the i-vectors onto a 400-dimensional subspace. The NFA-projected i-vectors
430 were then used to train PLDA models with 300 speaker factors ($Q_1 = 300$).
431 In our experiments, we make sure that the number of speaker factors plus
432 SNR and duration factors is no more than 400.

433 5.4. *SNR and Duration Groups*

434 To determine the SNR and duration subspaces in the SI-PLDA, DI-
435 PLDA and SDI-PLDA models, the development data described in Sec-
436 tion 5.1 were divided into K groups according to the measured SNRs and
437 duration of the utterances, where K varied from 3 to 8.⁴ Because SNR
438 and duration are continuous variables, there will be infinite possible ways of
439 dividing them into intervals. Therefore, we evenly divided the training utter-
440 ances into K groups such that each group contains almost the same number
441 of training i-vectors. Although this partitioning method leads to unequal
442 SNR and duration intervals, it ensures that each group has sufficient train-
443 ing i-vectors for estimating the SNR and duration loading matrices reliably.
444 Table 2 lists the SNR range and duration range when $K = 8$ and $P = 8$ in
445 Eq. 29.

⁴To be precise, the first $K - 1$ groups have $\lfloor \frac{N}{K} \rfloor$ i-vectors, whereas the last one contains $N - (K - 1) \lfloor \frac{N}{K} \rfloor$, where $\lfloor x \rfloor$ is the floor of x .

Table 2: Division of SNR and duration groups for the training data of female speakers in CC4, when $K = 8$ and $P = 8$ in Eq. 29.

Group	SNR Range (dB)	Duration Range (s)
1	2.0–5.5	3–70
2	5.5–7.4	70–88
3	7.4–10.5	88–103
4	10.5–14.6	103–117
5	14.6–19.8	117–133
6	19.8–34.4	133–152
7	34.4–39.1	152–184
8	39.1–55.0	184–1215

446 6. Results and Analysis

447 We used equal error rate (EER) and $\min C_{\text{primary}}$, which is the same as
 448 minimum normalized decision cost function ($\min \text{DCF}$) defined in NIST 2012
 449 SRE [27], to evaluate the performance of different PLDA models. Table 1
 450 summarizes their abbreviations and formulations.

451 6.1. Effectiveness of SNR and Duration Factors

452 The first experiment aims to compare the effectiveness of SI-PLDA and
 453 DI-PLDA in compensating SNR variability and duration variability, respec-
 454 tively. Table 3 shows the results of SI-PLDA and DI-PLDA on CC1 and
 455 CC4 for different numbers of SNR groups and duration groups. [The results](#)
 456 [show that both SI-PLDA and DI-PLDA outperform PLDA. This suggests](#)
 457 [that including the duration subspace in DI-PLDA and the SNR subspace in](#)
 458 [SI-PLDA enables these models to address mismatch caused by duration and](#)
 459 [SNR, respectively. Moreover, SI-PLDA not only outperforms DI-PLDA \(in](#)
 460 [EER\) in most cases, but also performs stably with respect to the number](#)
 461 [of groups \$K\$. On the other hand, the performance of DI-PLDA on CC1](#)
 462 [\(female\) drops when the number of groups increases to 8.](#)

463 The second experiment compares the proposed SDI-PLDA model with
 464 other PLDA models and PLDA with uncertainty propagation. Results in
 465 Table 3 show that SDI-PLDA achieves the best performance in terms of EER
 466 and $\min \text{DCF}$ on CC4. This result suggests that SDI-PLDA can compensate
 467 for SNR and duration variabilities in the i-vector space. While UP-PLDA
 468 achieves the best performance in CC1, it performs badly under CC4. The
 469 reason is that CC4 involves more test trials with long duration than CC1
 470 (as shown in Fig. 5). As the i-vectors corresponding to utterances of long

Table 3: Performance of PLDA, UP-PLDA, SI-PLDA, DI-PLDA and SDI-PLDA in CC1 and CC4 of NIST 2012 SRE core set. K and P denote the number of SNR and duration groups, respectively. The best results are highlighted in boldface.

Model	K	P	Male				Female			
			CC1		CC4		CC1		CC4	
			EER(%)	minDCF	EER(%)	minDCF	EER(%)	minDCF	EER(%)	minDCF
PLDA	–	–	5.28	0.374	2.69	0.317	7.39	0.514	2.35	0.332
UP-PLDA	–	–	3.89	0.346	3.55	0.493	5.47	0.408	3.36	0.483
SI-PLDA	3	–	5.28	0.369	2.56	0.292	7.06	0.504	2.18	0.299
	4	–	5.28	0.368	2.56	0.287	7.12	0.499	2.18	0.287
	5	–	5.15	0.395	2.50	0.288	7.07	0.498	2.12	0.291
	6	–	5.22	0.370	2.51	0.281	7.14	0.508	2.16	0.292
	7	–	5.36	0.381	2.48	0.284	7.13	0.505	2.16	0.287
	8	–	5.23	0.369	2.48	0.288	7.03	0.506	2.13	0.287
DI-PLDA	–	3	5.42	0.368	2.60	0.287	6.98	0.512	2.25	0.287
	–	4	5.57	0.369	2.55	0.291	6.98	0.499	2.23	0.299
	–	5	5.42	0.369	2.56	0.287	7.02	0.503	2.25	0.289
	–	6	5.21	0.369	2.52	0.289	7.30	0.526	2.16	0.293
	–	7	5.36	0.381	2.55	0.287	7.38	0.534	2.18	0.302
	–	8	5.23	0.369	2.56	0.289	8.95	0.563	2.20	0.288
SDI-PLDA	3	3	5.42	0.375	2.52	0.287	6.93	0.496	2.15	0.284
	4	4	5.41	0.368	2.54	0.289	6.97	0.491	2.13	0.288
	5	5	5.13	0.367	2.55	0.288	6.96	0.491	2.15	0.289
	6	6	5.14	0.367	2.49	0.283	7.05	0.508	2.14	0.289
	7	7	5.42	0.373	2.34	0.280	7.13	0.505	2.13	0.289
	8	8	5.54	0.373	2.49	0.286	8.08	0.556	2.11	0.284

471 duration are reliable, UP-PLDA loses its advantage in handling reliable i-
472 vectors compared to PLDA.

473 To confirm that SDI-PLDA really outperforms SI-PLDA and DI-PLDA,
474 we performed McNemar’s tests [9] on the differences between the EERs. For
475 each model, the best performing configuration (by varying K and P) was
476 used in the tests. The p-values of these tests are shown in Table 4. As the
477 p-values between SDI-PLDA and the other two models are less than 0.05,
478 we conclude that SDI-PLDA outperforms SI-PLDA and DI-PLDA.

479 We have also linearly fused the scores of the best performing DI-PLDA
480 and SI-PLDA in CC4, with fusion weights for DI-PLDA and SI-PLDA set
481 to 0.65 and 0.35, respectively. For male speakers, the EER after fusion is
482 2.41% and the minDCF is 0.282. For female speakers, the EER after fusion
483 is 2.13% and the minDCF is 0.283. This fusion performance is comparable

Table 4: P-values of McNemar’s tests [9] on the differences in EERs based on CC4 of NIST 2012 SRE core set, male speakers. For each entry, $p < 0.05$ means that the difference between the EERs is statistically significant at a confidence level of 95%.

Method	DI-PLDA	SDI-PLDA
SI-PLDA	0.002	0.013
DI-PLDA	–	0.022

Table 5: Performance of DI-PLDA on CC1 and CC4 of NIST 2012 SRE (male, core set) using different approaches to deriving the EM training algorithm and the scoring function. *EM*: EM is derived by assuming that the latent factors are posteriorly independent (Eqs. 14–24 of [21]). *VB-EM*: EM is derived by using variational Bayes and the latent factors are assumed posteriorly dependent (Eq. 6–Eq. 23). *Scoring1*: The duration is unknown during scoring (Eq. 26 in [21]). *Scoring2*: The duration is known during scoring, and the duration factor has sharp posterior (Eq. 26). *Scoring3*: The duration is known during scoring, and the duration factor has blunt posterior (Eq. 28). P in Eq. 4 was set to 6.

Method	CC1		CC4	
	EER(%)	minDCF	EER(%)	minDCF
EM + Scoring1	5.21	0.369	2.52	0.289
VB-EM + Scoring1	5.42	0.366	2.56	0.285
EM + Scoring2	5.42	0.371	2.58	0.290
EM + Scoring3	5.28	0.366	2.57	0.290

484 with that of SDI-PLDA, suggesting that SNR and duration variabilities can
 485 be handled either in the model domain (Eq. 29) or in the score domain. But
 486 the later requires a set of optimal fusion weights to achieve a performance
 487 comparable to that of the former.

488 6.2. Numbers of SNR and Duration Groups

489 Another observation from Table 3 is that the performance of DI-PLDA
 490 and SDI-PLDA on CC1 becomes worse when the numbers of SNR and du-
 491 ration groups, K and P , increase. This suggests that appropriate values of
 492 K and P are important for DI-PLDA and SDI-PLDA. Since the amount
 493 of training data for CC1 is much less than that for CC4, when K and
 494 P increase, the number of training samples in each group becomes limited,
 495 causing unreliable estimation of SNR and duration loading matrices. Hence,
 496 the values of K and P should be determined based on the amount of train-

Table 6: Performance of SDI-PLDA in CC4 of NIST 2012 SRE core set for male speakers with varying numbers of SNR and duration factors. The numbers of SNR and duration groups were fixed to 7.

Q_2 & Q_3	EER(%)	minDCF
10	2.34	0.280
20	2.33	0.281
30	2.35	0.283
40	2.37	0.286

ing data. In particular, if K and P are very large, there will be so many SNR factors and duration factors that each i-vector is considered to be obtained from a unique SNR or duration. This means that the SNR- and duration-invariant PLDA models reduce to the traditional Gaussian PLDA, which only considers the session variability instead of the variability caused by different SNRs and durations.

6.3. Combinations of Training and Scoring Methods

Table 5 shows the performance of DI-PLDA under different combinations of training methods and scoring methods. The results suggest that using the original training method (EM) and scoring method (Scoring1) in [21] achieves the best results, which assumes that the latent factors are posteriorly independent and that the SNR and duration of utterance are unknown. Although the EM algorithm derived from variational Bayes (VB) is more theoretically justifiable, VB-EM + Scoring1 in Table 5 does not outperform EM + Scoring1. This is a rather unexpected result. One possible reason is that because the number of training i-vectors in each duration group is large enough to make the posterior density of duration factors (\mathbf{y}_p in Eq. 4) very sharp, causing the joint posterior density $p(\mathbf{h}_i, \mathbf{y}_p | \mathbf{x})$ to spread mainly along \mathbf{h}_i instead of spreading over both \mathbf{h}_i and \mathbf{y}_p . As a result, the assumption that the latent factors \mathbf{h}_i and \mathbf{y}_p are posteriorly independent becomes valid.

Comparing Rows 1, 3, and 4 in Table 5 suggests that scoring with duration information does not achieve any advantage. This may be because the EM and VB-EM have already taken duration variability into account through the duration loading matrix.

6.4. Numbers of SNR and Duration Factors

To investigate the effect of varying the number of SNR and duration factors, we set Q_2 and Q_3 to different values but keeping K and P fixed.

Table 7: Performance (EER(%)/minDCF) of (a) PLDA and (b) SDI-PLDA on CC4 of NIST 2012 SRE (male core set) under different combinations of SNR (dB) and utterance durations. The last row shows the relative increases in EER/minDCF when SNR decreases from 15dB to 6dB. The last column shows the relative increases in EER/minDCF when number of frames reduces from 2500 to 1000. For each SNR and duration combination, the one with a smaller increase in EER or minDCF is highlighted in boldface.

SNR	Number of frames					Relative Inc.
	1000	1334	1667	2000	2500	
6	6.02/0.517	5.40/0.483	4.78/0.454	4.35/0.431	3.86/0.409	0.560/ 0.264
8	5.47/0.499	4.79/0.454	4.25/0.427	3.88/0.405	3.71/0.393	0.474/0.270
10	4.68/0.460	4.32/0.431	4.00/0.405	3.72/0.388	3.41/0.367	0.372/ 0.253
12	4.58/0.462	4.07/0.416	3.53/0.399	3.49/0.374	3.23/0.362	0.418/ 0.276
15	4.57/0.441	3.87/0.412	3.71/0.389	3.48/0.363	3.15/0.351	0.451/ 0.256
Relative Inc.	0.317/ 0.172	0.395/ 0.172	0.288/ 0.167	0.250/ 0.187	0.225/0.165	–

(a) PLDA

SNR	Number of frames					Relative Inc.
	1000	1334	1667	2000	2500	
6	5.53/0.473	5.03/0.437	4.39/0.410	4.03/0.398	3.66/0.360	0.512 /0.314
8	5.13/0.453	4.54/0.414	3.85/0.385	3.56/0.368	3.42/0.354	0.500/0.280
10	4.43/0.423	4.07/0.393	3.75/0.370	3.53/0.353	3.27/0.330	0.355 /0.282
12	4.29/0.421	3.77/0.379	3.39/0.358	3.28/0.330	3.17/0.314	0.353 /0.341
15	4.42/0.401	3.78/0.367	3.44/0.347	3.31/0.322	3.01/0.314	0.468/0.277
Relative Inc.	0.251 /0.179	0.331 /0.191	0.276 /0.182	0.218 /0.236	0.216 / 0.146	–

(b) SDI-PLDA

524 The effect is shown in Table 6. Although there is no obvious relation between
525 the number of speaker factors and the number of SNR factors and duration
526 factors, this table suggests that it is fine to set Q2 and Q3 to 10.

527 6.5. Robustness Against Mismatch Types

528 The results in Table 3 do not show which type of mismatches (SNR or
529 duration) is more harmful to performance. To this end, we fixed one type of
530 variability and vary the other type. The test utterances from male speakers
531 in CC4, excluding utterances with less than 2500 frames, were used as the
532 evaluation data. We added babble noise to the test utterances at SNR
533 of 6dB, 8dB, 10dB, 12dB, and 15dB. Then, 2500 frames were randomly
534 selected from each test utterance at each SNR. Finally, test utterances with

535 different durations were created by successively discarding some frames from
536 the 2500-frame test utterances. The SNRs and durations were set such that
537 the percentage decrease in successive SNR is equal to the percentage decrease
538 in successive utterance length. For example, when SNR reduces from 15dB
539 to 12dB, the number of frames decreases from 2500 to 2000, which amount
540 to 20% relative reduction.

541 Table 7 shows the results of PLDA and SDI-PLDA under different com-
542 binations of SNRs and utterance durations. The results show that system
543 performance degrades with decreasing SNR (from 15dB to 6dB) or utterance
544 duration (from 2500 frames to 1000 frames). Both tables show that there
545 is almost no change in performance once the SNR is larger than or equal to
546 12dB, which suggests that the effect of SNR variability is small when the
547 test utterances are not noisy.

548 Comparing the performance and the relative increases in EER and minDCF
549 between PLDA and SDI-PLDA in Table 7, we can draw the following conclu-
550 sions: (1) the SDI-PLDA performs better than PLDA for all combinations
551 of utterance length and SNR; and (2) when either one the two variabil-
552 ity types is fixed but the other is varied, the EER of SDI-PLDA is more
553 stable (smaller relative increase) but its minDCF is less stable (larger rela-
554 tive increase). Despite its larger relative increases in minDCF, SDI-PLDA
555 achieves a lower minDCF under all conditions, which provides strong evi-
556 dence supporting its superiority over PLDA in tackling SNR and duration
557 variabilities.

558 7. Conclusions

559 A new SNR- and duration-invariant PLDA model is presented. It is de-
560 signed to improve the robustness of speaker verification systems under both
561 noise-level and duration mismatches. By introducing a duration subspace to
562 SNR-invariant PLDA, duration information can be captured and the effect
563 of noise-level variability and duration variability can be simultaneously sup-
564 pressed. Experiments on the NIST 2012 SRE demonstrate the effectiveness
565 of the proposed method.

566 8. Acknowledgment

567 This work was supported in part by The RGC of Hong Kong SAR (Grant
568 Nos. PolyU 152117/14E and PolyU 152068/15E) and in part by the Taiwan
569 MOST with Grant 105-2221-E-009-137-MY2.

570 **Appendix A.**

571 To simplify notations, we use \mathbf{x}_s and \mathbf{x}_t instead of $\hat{\mathbf{x}}_s$ and $\hat{\mathbf{x}}_t$ in Eq. 51
 572 to represent the NFA-projected i-vectors. If \mathbf{x}_s and \mathbf{x}_t are from the same
 573 speaker, then we have

$$\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix} + \begin{bmatrix} \mathbf{V} & \mathbf{U} & \mathbf{0} & \mathbf{R} & \mathbf{0} \\ \mathbf{V} & \mathbf{0} & \mathbf{U} & \mathbf{0} & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{h} \\ \mathbf{w}_s \\ \mathbf{w}_t \\ \mathbf{y}_s \\ \mathbf{y}_t \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_s \\ \boldsymbol{\epsilon}_t \end{bmatrix}, \quad (\text{A.1})$$

574 where \mathbf{h} represents the speaker factor shared by both i-vectors, \mathbf{w}_s and \mathbf{w}_t
 575 represent the SNR factors of the two utterances, and \mathbf{y}_s and \mathbf{y}_t represent the
 576 duration factors of the two utterances, respectively. Eq. A.1 can be written
 577 in a compact form:

$$\tilde{\mathbf{x}}_{st} = \tilde{\mathbf{m}} + \tilde{\mathbf{A}}\tilde{\mathbf{z}}_{st} + \tilde{\boldsymbol{\epsilon}}_{st}$$

578 where the tilde denotes the stacking of vectors and

$$\tilde{\mathbf{A}} = \begin{bmatrix} \mathbf{V} & \mathbf{U} & \mathbf{0} & \mathbf{R} & \mathbf{0} \\ \mathbf{V} & \mathbf{0} & \mathbf{U} & \mathbf{0} & \mathbf{R} \end{bmatrix}.$$

Assuming that the NFA-projected i-vectors follow a Gaussian distribution, the distribution of $\tilde{\mathbf{x}}_{st}$ can be obtained by marginalizing over all possible latent factors as follows:

$$\begin{aligned} p(\tilde{\mathbf{x}}_{st}|\text{same-speaker}) &= \int p(\tilde{\mathbf{x}}_{st}|\tilde{\mathbf{z}}_{st})p(\tilde{\mathbf{z}}_{st})d\tilde{\mathbf{z}}_{st} \\ &= \int \mathcal{N}(\tilde{\mathbf{x}}_{st}|\tilde{\mathbf{m}} + \tilde{\mathbf{A}}\tilde{\mathbf{z}}_{st}, \tilde{\boldsymbol{\Sigma}})\mathcal{N}(\tilde{\mathbf{z}}_{st}|\mathbf{0}, \mathbf{I})d\tilde{\mathbf{z}}_{st} \\ &= \mathcal{N}(\tilde{\mathbf{x}}_{st}|\tilde{\mathbf{m}}, \tilde{\mathbf{A}}\tilde{\mathbf{A}}^\top + \tilde{\boldsymbol{\Sigma}}) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{tot} & \boldsymbol{\Sigma}_{ac} \\ \boldsymbol{\Sigma}_{ac} & \boldsymbol{\Sigma}_{tot} \end{bmatrix}\right) \end{aligned} \quad (\text{A.2})$$

579 where $\tilde{\Sigma} = \text{diag}\{\Sigma, \Sigma\}$, $\Sigma_{tot} = \mathbf{V}\mathbf{V}^\top + \mathbf{U}\mathbf{U}^\top + \mathbf{R}\mathbf{R}^\top + \Sigma$ and $\Sigma_{ac} = \mathbf{V}\mathbf{V}^\top$.
 580 If \mathbf{x}_s and \mathbf{x}_t are from the utterances of two different speakers, we have

$$\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} = \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix} + \begin{bmatrix} \mathbf{V} & \mathbf{0} & \mathbf{U} & \mathbf{0} & \mathbf{R} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} & \mathbf{0} & \mathbf{U} & \mathbf{0} & \mathbf{R} \end{bmatrix} \begin{bmatrix} \mathbf{h}_s \\ \mathbf{h}_t \\ \mathbf{w}_s \\ \mathbf{w}_t \\ \mathbf{y}_s \\ \mathbf{y}_t \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_s \\ \boldsymbol{\epsilon}_t \end{bmatrix} \quad (\text{A.3})$$

581 which can be compactly written as

$$\tilde{\mathbf{x}}_{st} = \tilde{\mathbf{m}} + \bar{\mathbf{A}}\bar{\mathbf{z}}_{st} + \tilde{\boldsymbol{\epsilon}}_{st}.$$

The distribution of $\tilde{\mathbf{x}}_{st}$ is obtained by marginalizing over $\bar{\mathbf{z}}_{st}$:

$$\begin{aligned} p(\tilde{\mathbf{x}}_{st} | \text{diff-speaker}) &= \int p(\tilde{\mathbf{x}}_{st} | \bar{\mathbf{z}}_{st}) p(\bar{\mathbf{z}}_{st}) d\bar{\mathbf{z}}_{st} \\ &= \int \mathcal{N}(\tilde{\mathbf{x}}_{st} | \tilde{\mathbf{m}} + \bar{\mathbf{A}}\bar{\mathbf{z}}_{st}, \tilde{\Sigma}) \mathcal{N}(\bar{\mathbf{z}}_{st} | \mathbf{0}, \mathbf{I}) d\bar{\mathbf{z}}_{st} \\ &= \mathcal{N}(\tilde{\mathbf{x}}_{st} | \tilde{\mathbf{m}}, \bar{\mathbf{A}}\bar{\mathbf{A}}^\top + \tilde{\Sigma}) \\ &= \mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} \Sigma_{tot} & \mathbf{0} \\ \mathbf{0} & \Sigma_{tot} \end{bmatrix}\right). \end{aligned} \quad (\text{A.4})$$

Combining Eq. A.2 and Eq. A.4, we have the log-likelihood ratio score:

$$\begin{aligned} S_{\text{LR}}(\mathbf{x}_s, \mathbf{x}_t) &= \ln \frac{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} \Sigma_{tot} & \Sigma_{ac} \\ \Sigma_{ac} & \Sigma_{tot} \end{bmatrix}\right)}{\mathcal{N}\left(\begin{bmatrix} \mathbf{x}_s \\ \mathbf{x}_t \end{bmatrix} \middle| \begin{bmatrix} \mathbf{m} \\ \mathbf{m} \end{bmatrix}, \begin{bmatrix} \Sigma_{tot} & \mathbf{0} \\ \mathbf{0} & \Sigma_{tot} \end{bmatrix}\right)} \\ &= \frac{1}{2} [\bar{\mathbf{x}}_s^\top \quad \bar{\mathbf{x}}_t^\top] \begin{bmatrix} \mathbf{Q} & \mathbf{P} \\ \mathbf{P} & \mathbf{Q} \end{bmatrix} \begin{bmatrix} \bar{\mathbf{x}}_s \\ \bar{\mathbf{x}}_t \end{bmatrix} + \text{const} \\ &= \frac{1}{2} [\bar{\mathbf{x}}_s^\top \mathbf{Q} \bar{\mathbf{x}}_s + 2\bar{\mathbf{x}}_s^\top \mathbf{P} \bar{\mathbf{x}}_t + \bar{\mathbf{x}}_t^\top \mathbf{Q} \bar{\mathbf{x}}_t] + \text{const}, \end{aligned} \quad (\text{A.5})$$

582 where

$$\begin{aligned} \bar{\mathbf{x}}_s &= \mathbf{x}_s - \mathbf{m}, & \bar{\mathbf{x}}_t &= \mathbf{x}_t - \mathbf{m}, \\ \mathbf{P} &= \Sigma_{tot}^{-1} \Sigma_{ac} (\Sigma_{tot} - \Sigma_{ac} \Sigma_{tot}^{-1} \Sigma_{ac})^{-1}, \\ \mathbf{Q} &= \Sigma_{tot}^{-1} - (\Sigma_{tot} - \Sigma_{ac} \Sigma_{tot}^{-1} \Sigma_{ac})^{-1}. \end{aligned}$$

583 **References**

- 584 [1] Billingsley, P., 2008. Probability and Measure. John Wiley & Sons,
585 New York.
- 586 [2] Bishop, C., 2006. Pattern Recognition and Machine Learning. Springer,
587 New York.
- 588 [3] Cai, W., Li, M., Li, L., Hong, Q., 2015. Duration dependent covariance
589 regularization in PLDA modeling for speaker verification, in: Proc.
590 Interspeech, pp. 1027–1031.
- 591 [4] Cumani, S., 2015. Fast scoring of full posterior PLDA models.
592 IEEE/ACM Transactions on Audio, Speech, and Language Processing
593 23, 2036–2045.
- 594 [5] Cumani, S., Plchot, O., Laface, P., 2014. On the use of i-vector
595 posterior distributions in probabilistic linear discriminant analysis.
596 IEEE/ACM Transactions on Audio, Speech, and Language Processing
597 22, 846–857.
- 598 [6] Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., Ouellet, P., 2011.
599 Front-end factor analysis for speaker verification. IEEE Trans. on Au-
600 dio, Speech, and Language Processing 19, 788–798.
- 601 [7] Garcia-Romero, D., Espy-Wilson, C., 2011. Analysis of i-vector length
602 normalization in speaker recognition systems, in: Proc. Interspeech, pp.
603 249–252.
- 604 [8] Garcia-Romero, D., Zhou, X., Espy-Wilson, C., 2012. Multicondition
605 training of Gaussian PLDA models in i-vector space for noise and rever-
606 beration robust speaker recognition, in: Proc. ICASSP, pp. 4257–4260.
- 607 [9] Gillick, L., Cox, S.J., 1989. Some statistical issues in the comparison
608 of speech recognition algorithms, in: Proc. ICASSP, Glasgow, UK. pp.
609 532–535.
- 610 [10] Hasan, T., Sadjadi, S.O., Liu, G., Shokouhi, N., Boril, H., Hansen, J.H.,
611 2013a. CRSS systems for 2012 NIST speaker recognition evaluation, in:
612 Proc. ICASSP, pp. 6783–6787.
- 613 [11] Hasan, T., Saeidi, R., Hansen, J.H., van Leeuwen, D., et al., 2013b.
614 Duration mismatch compensation for i-vector based speaker recognition
615 systems, in: Proc. ICASSP, pp. 7663–7667.

- 616 [12] Hatch, A., Kajarekar, S., Stolcke, A., 2006. Within-class covariance
617 normalization for SVM-based speaker recognition, in: Proc. ICSLP,
618 pp. 1471–1474.
- 619 [13] Hong, Q., Li, L., Li, M., Huang, L., Wan, L., Zhang, J., 2015. Modified-
620 prior PLDA and score calibration for duration mismatch compensation
621 in speaker recognition system, in: Proc. Interspeech, pp. 1037–1041.
- 622 [14] Kanagasundaram, A., Dean, D., Sridharan, S., Gonzalez-Dominguez,
623 J., Gonzalez-Rodriguez, J., Ramos, D., 2014. Improving short utter-
624 ance i-vector speaker verification using utterance variance modelling
625 and compensation techniques. *Speech Communication* 59, 69–82.
- 626 [15] Kenny, P., 2005. Joint factor analysis of speaker and session variability:
627 Theory and algorithms. CRIM, Montreal,(Report) CRIM-06/08-13 .
- 628 [16] Kenny, P., 2010. Bayesian speaker verification with heavy-tailed priors,
629 in: Proc. of Odyssey: Speaker and Language Recognition Workshop,
630 Brno, Czech Republic.
- 631 [17] Kenny, P., Boulianne, G., Ouellet, P., Dumouchel, P., 2006. Improve-
632 ments in factor analysis based speaker verification, in: Proc. ICASSP,
633 pp. 113–116.
- 634 [18] Kenny, P., Ouellet, P., Dehak, N., Gupta, V., Dumouchel, P., 2008. A
635 study of inter-speaker variability in speaker verification. *IEEE Trans.
636 on Audio, Speech and Language Processing* 16, 980–988.
- 637 [19] Kenny, P., Stafylakis, T., Ouellet, P., Alam, M.J., Dumouchel, P., 2013.
638 PLDA for speaker verification with utterances of arbitrary duration, in:
639 Proc. ICASSP, pp. 7649–7653.
- 640 [20] Li, N., Mak, M.W., 2015a. SNR-invariant PLDA modeling for robust
641 speaker verification, in: Proc. Interspeech, pp. 2317–2321.
- 642 [21] Li, N., Mak, M.W., 2015b. SNR-invariant PLDA modeling in nonpara-
643 metric subspace for robust speaker verification. *IEEE/ACM Trans. on
644 Audio, Speech and Language Processing* 23, 1648–1659.
- 645 [22] Li, Z., Lin, D., Tang, X., 2009. Nonparametric discriminant analysis for
646 face recognition. *IEEE Transactions on Pattern Analysis and Machine
647 Intelligence* 31, 755–761.

- 648 [23] Mak, M.W., 2014. SNR-dependent mixture of PLDA for noise robust
649 speaker verification, in: Proc. Interspeech, pp. 1855–1859.
- 650 [24] Mak, M.W., Pang, X.M., Chien, J.T., 2016. Mixture of PLDA for
651 noise robust i-vector speaker verification. *IEEE/ACM Trans. on Audio,
652 Speech and Language Processing* 24, 130–142.
- 653 [25] Mak, M.W., Yu, H.B., 2014. A study of voice activity detection tech-
654 niques for NIST speaker recognition evaluations. *Computer Speech and
655 Language* 28, 295–313.
- 656 [26] McLaren, M., Mandasari, M., Leeuwen, D., 2012. Source normalization
657 for language-independent speaker recognition using i-vectors, in: Proc.
658 Odyssey, pp. 55–61.
- 659 [27] NIST, 2012. The NIST year 2012 speaker recognition evaluation plan.
660 <http://www.nist.gov/itl/iad/mig/sre12.cfm> .
- 661 [28] Norwich, K.H., 1993. Information, sensation, and perception. Academic
662 Press San Diego.
- 663 [29] Pelecanos, J., Sridharan, S., 2001. Feature warping for robust speaker
664 verification, in: Proc. Odyssey: The Speaker and Language Recognition
665 Workshop, Crete, Greece. pp. 213–218.
- 666 [30] Prince, S.J., 2012. *Computer Vision: Models, Learning, and Inference*.
667 Cambridge University Press.
- 668 [31] Prince, S.J., Elder, J.H., 2007. Probabilistic linear discriminant anal-
669 ysis for inferences about identity, in: Proc. IEEE 11th International
670 Conference on Computer Vision, IEEE. pp. 1–8.
- 671 [32] Sarkar, A.K., Matrouf, D., Bousquet, P.M., Bonastre, J.F., 2012. Study
672 of the effect of i-vector modeling on short and mismatch utterance du-
673 ration for speaker verification., in: Proc. Interspeech, pp. 2662–2665.
- 674 [33] Sizov, A., Lee, K.A., Kinnunen, T., 2014. Unifying probabilistic linear
675 discriminant analysis variants in biometric authentication, in: *Struc-
676 tural, Syntactic, and Statistical Pattern Recognition*. Springer, pp. 464–
677 475.
- 678 [34] Vesnicer, B., Zganec-Gros, J., Dobrisek, S., Struc, V., 2014. Incorporating
679 duration information into i-vector-based speaker-recognition
680 systems, in: Proc. Odyssey, pp. 241–248.

- 681 [35] Watanabe, S., Chien, J.T., 2015. Bayesian Speech and Language Pro-
682 cessing. Cambridge University Press.
- 683 [36] Yamamoto, H., Koshinaka, T., 2015. Denoising autoencoder-based
684 speaker feature restoration for utterances of short duration, in: Proc.
685 Interspeech, pp. 1052–1056.
- 686 [37] Yu, H., Mak, M., 2011. Comparison of voice activity detectors for
687 interview speech in NIST speaker recognition evaluation, in: Proc. In-
688 terspeech, pp. 2353–2356.