
Variational Autoencoders with implicit priors for short-duration text-independent speaker verification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this work, we exploited different strategies to provide prior knowledge to com-
2 monly used generative modeling approaches aiming to obtain speaker-dependent
3 low dimensional representations from short-duration segments of speech data,
4 making use of available information of speaker identities. Namely, convolutional
5 variational autoencoders are employed, and statistics of its learned posterior distri-
6 bution are used as low dimensional representations of fixed length short-duration
7 utterances. In order to enforce speaker dependency in the latent layer, we intro-
8 duced a variation of the commonly used prior within the variational autoencoders
9 framework, i.e. the model is simultaneously trained for reconstruction of inputs
10 along with a discriminative task performed on top of latent layers outputs. The
11 effectiveness of both triplet loss minimization and speaker recognition are evaluated
12 as implicit priors on the challenging cross-language NIST SRE 2016 setting and
13 compared against fully supervised and unsupervised baselines.

14 1 Introduction

15 Variational autoencoders (VAEs) (1; 2) have been introduced as an effective framework within the
16 context of generative models that support tractable approximate inference (3), leveraging neural
17 networks both for generative modeling as well as for approximate inference, usually employing a
18 non-informative prior. However, follow-up works have shown that too simplistic of a prior will in
19 general lead to also simplistic posteriors which might not encode relevant information about the
20 inputs.

21 Attempts to overcome the above mentioned issue include adversarial autoencoders, proposed origi-
22 nally in (4), which employ an adversarial game on top of latent variables. The discriminator tries
23 to distinguish samples from the posterior and prior distributions, while the encoder tries to produce
24 samples that are indistinguishable from the prior. Moreover, stochastic variational methods (5; 6)
25 appeared as an alternative in which informative data-dependent priors can be used. Sampling methods
26 are employed to estimate gradients of the variational gap, such that any prior from which one can
27 sample can be used. In both of the described cases, the only requirement for a prior is the possibility
28 of efficiently sampling from it.

29 Even though aforementioned adversarial autoencoders and stochastic variational methods allow the
30 use of non-trivial priors, designing prior distributions which yield desired properties on the finally
31 learned posterior is a challenging task in itself. In this work, rather than explicitly matching posterior
32 and prior distributions, we evaluate the effectiveness of enforcing relevant properties on the posterior
33 distribution by introducing auxiliary discriminative tasks at train time, making use of available labels.
34 By doing so, we argue prior knowledge is introduced implicitly, since desired properties are directly
35 enforced into the posterior distribution.

36 The remainder of this paper is organized as follows: Section 2 includes a brief description of the VAE
 37 framework along with a brief definition of the speaker verification problem, which we employ as a
 38 test-bed for the proposed approach. Section 3 details the strategy we proposed in order to enforce
 39 desired properties within the VAE’s learned posterior. In Section 4 we evaluate our method, and
 40 finally draw conclusions along with future directions in Section 5.

41 2 Background: Variational Autoencoders and Speaker Verification

42 Consider $p(X, Z)$, where X is the observed data and Z is the latent representation. The posterior
 43 distribution $p(Z|X)$ can be approximated within the family of distributions $q(Z|X, \lambda)$, parametrized
 44 by λ . The so-called variational gap has to be minimized in order to give the maximum likelihood
 45 estimate of λ . The variational gap is defined as the Kullback-Leibler divergence between the
 46 approximate $q(Z|X, \lambda)$ and the true posterior over Z , $p(Z|X)$, written as $\text{KL}(q(Z|X, \lambda)||p(Z|X))$.

47 A common approach to minimize $\text{KL}(q_\lambda(Z|X)||p(Z|X))$ with respect to λ is to define the Evidence
 48 Lower Bound (ELBO) given by:

$$\text{ELBO}(\lambda) = \log(p(X)) - \text{KL}(q(Z|X, \lambda)||p(Z|X)), \quad (1)$$

49 whose terms can be rearranged, and ELBO can be simplified to:

$$\text{ELBO}(\lambda) = \mathbb{E}_q[\log p(X|Z)] - \text{KL}(\log q(Z|X, \lambda)||p(Z)). \quad (2)$$

50 Two main components present in above equation are the inference model $q(Z|X, \lambda)$ and the generative
 51 model $p(X|Z)$. VAEs parametrize both distributions using neural networks in an encoder/decoder
 52 setup. The encoder takes samples from X and outputs the parameters λ of the latent variable model
 53 $q_\theta(Z|X)$. The decoder receives samples from Z as input and returns reconstructed data samples from
 54 $p_\phi(X|Z)$. Parameters θ and ϕ are the weights and biases of the neural networks which are selected
 55 to minimize the negative ELBO using stochastic gradient descent. The negative of the ELBO yields
 56 the following loss function used for training the neural networks:

$$l(\theta, \phi) = -\mathbb{E}_{q_\theta(z|x)}[\log p_\phi(X|Z)] + \text{KL}(\log q_\theta(Z|X)||p(Z)). \quad (3)$$

57 First term in above equation is equivalent to maximum likelihood estimation, thus it is in general
 58 substituted by a reconstruction loss, while the second term can be seen as a regularizer, which tries to
 59 ensure that the approximation follows the prior distribution as much as possible.

60 The posterior $q_\theta(Z|X)$ is in general assumed to be an uncorrelated Gaussian. In order to train the VAE
 61 using stochastic gradient descent, the reparametrization trick (7; 8) is employed allowing gradients
 62 computation through the sampling process between encoder and decoder. Hence, the outputs of the
 63 encoder network are the statistics of $q_\theta(Z|X)$ and Z - input for the decoder - is ultimately obtained
 64 by $Z = \mu(X) + \sigma(X) \cdot \epsilon$, where $\mu(X)$ and $\sigma(X)$ are the encoder’s outputs given X , while ϵ is
 65 sampled from $\mathcal{N}(0, I)$.

66 Speaker verification consists of accepting or rejecting a claimed identity by comparing two spoken
 67 utterances, the first of these utterances being used for enrollment (produced by the speaker with the
 68 target identity) and the second utterance is obtained from the verified speaker (9).

69 Under the text-independent setting, speaker verification is performed on top of unconstrained spoken
 70 phrases of arbitrary length. The added phonetic variability in this scenario represents an extra adverse
 71 factor when compared to the session and speaker variabilities, present in the text-dependent case
 72 (10). Classical approaches for Automatic Speaker Verification split the problem into two distinct
 73 phases: (i) compute low dimensional speaker representations; (ii) perform binary classification on
 74 top of pre-computed representation of enrollment and test utterances.

75 3 Proposed Model

76 Unlike the ELBO-based loss definition in Equation 3, we evaluate the use of an auxiliary task on
77 top of the posterior statistics $\mu(X)$, with the aim at enforcing a multi-modal posterior with modes
78 depending on given class labels. Our training loss is thus defined by:

$$l(\theta, \phi) = (1 - \beta)\|X - X'\|_2^2 + \beta D(\mu(X), y), \quad (4)$$

79 where the first term, the mean squared error between the input X and its reconstructed pair X' , is
80 the same as in the standard VAE setting, while the second term, $D(\mu(X), y)$, is some discriminative
81 loss which plays the role of the KL term in Equation 3, considering given class labels y . $\beta \in [0, 1]$
82 is a tunable hyperparameter. $\mu(X)$ is employed as a low-dimensional embedding of inputs for the
83 discriminative auxiliary task. Two distinct choices of $D(\mu(X), y)$ are evaluated here:

- 84 1. A soft triplet loss defined by $\text{softplus}(\|d_+ - d_-\|_2)$, where d_+ and d_- correspond to
85 a distance measure between pairs of embeddings. d is chosen as $d(\mu(X_1), \mu(X_2)) =$
86 $1 - \frac{\mu(X_1) \cdot \mu(X_2)}{\|\mu(X_1)\|_2 \|\mu(X_2)\|_2}$, and the second term is the cosine of the smallest angle between
87 $\mu(X_1)$ and $\mu(X_2)$.
- 88 2. The sum of triplet loss with a multi-class classification loss, i.e. $\mu(X)$ is linearly projected
89 into an output layer and cross-entropy loss is measured using available labels.

90 We evaluate the described setting on the speaker verification task. *RMSProp* is employed for
91 optimization with α set to 0.99. The global learning rate starts at 0.001 and is halved once triplet
92 loss, measured on a validation set held out of training, plateaus for 30 epochs. Training is carried
93 out in a single Titan X NVIDIA GPU, with minibatches of size 64. Minibatches are constructed such
94 that two random segments of different utterances belonging to the same speaker are sampled to form
95 same class pairs (positive), and a random sample from a different speaker is selected to compose the
96 different classes pair (negative). β was at 0.8 for all experiments.

97 4 Experimental Setup and Results

98 Evaluation is performed on top of the cross-language NIST SRE 2016 setting (11). Test data in
99 Tagalog and Cantonese are available, while train data is in English. Embeddings obtained with a
100 standard VAE, along with our two proposed strategies using two distinct $D(\mu(X), y)$ previously
101 described choices are compared with x-vectors, a fully-supervised approach shown to outperform
102 i-vectors (12) in the full-recording setting (13). Train data is composed of: *Switchboard-2*, phases
103 1, 2, and 3, along with *NIST SREs* from 2004 to 2010 combined with *Mixer 6*, which sums up to
104 approximately 7000 speakers, out of which we remove all the recordings of 50 speakers to be used as
105 validation set. Training is performed on top of 40-dimensional log-mel filter banks. Only the SRE
106 portion is used for training probabilistic linear discriminant analysis (PLDA) (14), which was used as
107 a backend at evaluation phase. Since our model requires fixed size inputs, speech segments of 256
108 frames were randomly selected from each recording at train time. We augment the described train
109 dataset following the approach in (13), i.e. with additive background noise from the MUSAN corpus
110 and reverberation by convolving room impulse responses (RIR) with original audio data (MUSAN
111 and RIR are available at www.openslr.org). We removed silence frames from data using a simple
112 energy-based voice activity detector.

113 For enrollment, test, and unlabelled (used for PLDA adaptation) data, embeddings of each recording
114 are obtained from 256 frames windows without overlap, and then averaged, such that each test
115 utterance is represented by a single fixed dimensional representation, even though models only have
116 access to short-duration segments.

117 PLDA was employed as backend after dimensionality reduction of embeddings from 256 to 128,
118 using linear discriminant analysis. PLDA is trained on embeddings from the *SRE* partition of training
119 data, which are computed following the same approach as described for test data for the case of
120 our proposed models, while using the full-recordings in the case of x-vectors. Results in terms of
121 Equal Error Rate (EER) are shown in Table 1 for embeddings obtained from VAEs trained both in a
122 standard fashion, and our proposed approaches.

Table 1: EER obtained for embeddings averaged over short short-duration segments.

	PLDA			Adapted PLDA		
	Cantonese	Tagalog	Pooled	Cantonese	Tagalog	Pooled
X-vector	30.91	31.32	31.04	14.41	20.98	17.62
VAE	31.55	32.13	31.83	31.10	32.24	31.66
VAE+Triplet loss	21.81	27.80	24.79	19.89	25.50	22.76
VAE+Cross-entropy	21.46	27.05	24.28	16.50	23.00	20.02

123 As expected, including speaker identities relevantly increases the discriminability of learned repre-
 124 sentations when compared to a fully-unsupervised VAE, in both Tagalog and Cantonese evaluations.
 125 We further notice that performing speaker recognition on top of statistics of the posterior is more
 126 effective than the metric learning approach of triplet loss minimization alone.

127 In order to overcome the relevant domain shift between train and test data due to different spoken
 128 languages, the model adaptation scheme introduced in (15) is utilized for PLDA. To do so, embeddings
 129 of unlabelled data in Cantonese and Tagalog are clustered, and clusters are used as speaker identities,
 130 which are then employed for training a second PLDA model. The final model is obtained by simply
 131 averaging the second order statistics of the two trained models.

132 Results, as reported in the right section of Table 1, correspond to the evaluation using the adapted
 133 PLDA model. Interestingly, one can notice that the higher the *level of supervision* employed on
 134 embeddings model training, the higher is the performance gain when adaptation is used. By *level of*
 135 *supervision* we mean how relevant class labels (speaker identities in the studied case) are at train time.
 136 Standard VAE makes no use of class labels, while triplet loss employs such information for triplets
 137 construction only. Even in the case in which our VAE is trained with cross-entropy minimization,
 138 semi-supervised settings can be used, leveraging available unlabelled data, which is not the case for
 139 x-vectors, for instance, whose training is performed in a fully-supervised fashion. We thus argue that
 140 an increasing level of supervision induces domain-dependent representations, and this is the reason
 141 adaptation yields a huge improvement in such cases.

142 We further evaluate the discriminability of the representations corresponding to the statistics of
 143 posterior distributions approximated by VAEs trained in a standard fashion and making use of
 144 available speaker identities by plotting 2-dimensional t-SNE embeddings of $\mu(X)$, computed for 10
 145 speakers held out of training. Figures 1, 2, 3 are ordered in increasing *level of supervision*, which
 146 once more supports the claim that making use of class labels to perform discriminative tasks on top
 147 of statistics of the posterior is an effective strategy to enforce desired properties.

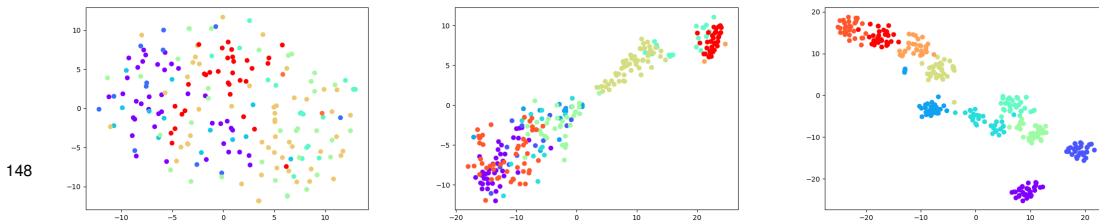


Figure 1: Embeddings obtained from a standard VAE posterior.

Figure 2: Embeddings obtained from a VAE trained with triplet loss minimization.

Figure 3: Embeddings from a VAE trained with cross-entropy minimization

149 5 Conclusion

150 In this work, we proposed to exchange the divergence term within the variational autoencoders
 151 training loss by some discriminative cost, leveraging available class labels. We thus argue such an
 152 approach is equivalent to implicitly defining prior distributions, directly inducing desired properties
 153 in the learned posterior distribution. Evaluation is performed on the challenging cross-language NIST
 154 SRE 2016 evaluation setting, for which we show embeddings obtained by such an approach are
 155 speaker-dependent, as enforced by discriminative tasks performed at train time. Future directions
 156 include the evaluation of this framework on the semi-supervised setting, employing unlabelled data
 157 for training of the generative model, along with labelled data.

References

- 158
- 159 [1] Diederik P Kingma and Max Welling, “Stochastic gradient vb and the variational auto-encoder,”
160 in *Second International Conference on Learning Representations, ICLR*, 2014.
- 161 [2] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra, “Stochastic backpropagation
162 and approximate inference in deep generative models,” *arXiv preprint arXiv:1401.4082*, 2014.
- 163 [3] Alex Lamb, Vincent Dumoulin, and Aaron Courville, “Discriminative regularization for
164 generative models,” *arXiv preprint arXiv:1602.03220*, 2016.
- 165 [4] Alireza Makhzani, Jonathon Shlens, Navdeep Jaitly, Ian Goodfellow, and Brendan Frey, “Ad-
166 versarial autoencoders,” *arXiv preprint arXiv:1511.05644*, 2015.
- 167 [5] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley, “Stochastic variational
168 inference,” *The Journal of Machine Learning Research*, vol. 14, no. 1, pp. 1303–1347, 2013.
- 169 [6] Rajesh Ranganath, Sean Gerrish, and David Blei, “Black box variational inference,” in *Artificial
170 Intelligence and Statistics*, 2014, pp. 814–822.
- 171 [7] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” *arXiv
172 preprint arXiv:1412.6980*, 2014.
- 173 [8] Danilo Jimenez Rezende and Shakir Mohamed, “Variational inference with normalizing flows,”
174 *arXiv preprint arXiv:1505.05770*, 2015.
- 175 [9] Wei Li, Tianfan Fu, Hanxu You, Jie Zhu, and Ning Chen, “Feature sparsity analysis for i-vector
176 based speaker verification,” *Speech Communication*, vol. 80, pp. 60–70, 2016.
- 177 [10] Tomi Kinnunen and Haizhou Li, “An overview of text-independent speaker recognition: From
178 features to supervectors,” *Speech communication*, vol. 52, no. 1, pp. 12–40, 2010.
- 179 [11] “NIST speaker recognition evaluation 2016,” 2016, [https://www.nist.gov/itl/iad/mig/speaker-
180 recognition-evaluation2016/](https://www.nist.gov/itl/iad/mig/speaker-recognition-evaluation2016/).
- 181 [12] Najim Dehak, Patrick J Kenny, Réda Dehak, Pierre Dumouchel, and Pierre Ouellet, “Front-end
182 factor analysis for speaker verification,” *IEEE Transactions on Audio, Speech, and Language
183 Processing*, vol. 19, no. 4, pp. 788–798, 2011.
- 184 [13] D. Snyder, D. Garcia-Romero, G. Sell, D. Povey, and S. Khudanpur, “X-vectors: Robust dnn
185 embeddings for speaker recognition,” in *2018 IEEE International Conference on Acoustics,
186 Speech and Signal Processing (ICASSP)*. IEEE, 2018.
- 187 [14] Simon JD Prince and James H Elder, “Probabilistic linear discriminant analysis for inferences
188 about identity,” in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*.
189 IEEE, 2007, pp. 1–8.
- 190 [15] Daniel Garcia-Romero, Alan McCree, Stephen Shum, Niko Brummer, and Carlos Vaquero,
191 “Unsupervised domain adaptation for i-vector speaker recognition,” in *Proceedings of Odyssey:
192 The Speaker and Language Recognition Workshop*, 2014.

193 **Appendix A - Model architecture**

194 Architectures employed for encoder and decoder are detailed in Tables 2 and 3. Batch normalization
 195 is used after all convolution layers. Inputs present dimensionality [40, 256], corresponding to 40 filter
 196 banks and 256 frames.

Table 2: Encoder architecture

Layer	Outputs	Kernel size	Stride	Dilation	Activation
<i>Convolution</i>	19, 84, 128	5, 5	2, 3	1, 2	ELU
<i>Convolution</i>	9, 40, 256	5, 5	2, 2	1, 2	ELU
<i>Convolution</i>	4, 40, 512	5, 5	2, 1	1, 1	ELU
<i>Convolution</i>	1, 40, 1024	5, 5	2, 1	1, 1	ELU
Average Pooling	1, 1, 1024	1, 40	1, 1	-	-
<i>Dense</i>	1024	-	-	-	ELU
<i>Dense</i>	256, 256	-	-	-	ELU, -

Table 3: Decoder architecture

Layer	Outputs	Kernel size	Stride	Dilation	Activation
<i>Dense</i>	800	-	-	-	ELU
<i>Transpose Convolution</i>	7, 14, 128	3, 4	1, 2	1, 3	ELU
<i>Transpose Convolution</i>	11, 29, 128	3, 4	2, 2	1, 2	ELU
<i>Transpose Convolution</i>	19, 59, 256	3, 4	2, 2	1, 2	ELU
<i>Transpose Convolution</i>	18, 118, 128	4, 6	1, 2	1, 1	ELU
<i>Transpose Convolution</i>	38, 246, 32	4, 12	2, 2	1, 1	ELU
<i>Transpose Convolution</i>	40, 256, 1	5, 13	1, 1	1, 1	-