



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

ScienceDirect

International Journal of Impact Engineering xxx (2016) xxx-xxx

[www.elsevier.com/locate/ijimpeng](http://www.elsevier.com/locate/ijimpeng)

# Sparse coding over redundant dictionaries for fast adaptation of speech recognition system

S. Shahnawazuddin\*, Rohit Sinha

*Department of Electronics and Electrical Engineering, Indian Institute of Technology Guwahati, Guwahati 781039, India*

Received 23 April 2015; received in revised form 4 April 2016; accepted 28 October 2016

## Abstract

This work presents a novel use of the sparse coding over redundant dictionary for fast adaptation of the acoustic models in the hidden Markov model-based automatic speech recognition (ASR) systems. The presented work is an extension of the existing acoustic model-interpolation-based fast adaptation approaches. In these methods, the basis (model) weights are estimated using an iterative procedure employing the maximum-likelihood (ML) criterion. For effective adaptation, typically a number of bases are selected and as a result of that the latency of the iterative weight estimation process becomes high for those ASR tasks that involve human-machine interactions. To address this issue, we propose the use of sparse coding of the target mean supervector over a speaker-specific (exemplar) redundant dictionary. In this approach, the employed greedy sparse coding not only selects the desired bases but also compresses them into a single supervector, which is then ML scaled to yield the adapted mean parameters. Thus reducing the latency in the basis weight estimation in comparison to the existing fast adaptation techniques. Further, to address the loss in information due to reduced degrees of freedom, we have also extended the proposed approach using separate sparse codings over multiple (exemplar and learned) redundant dictionaries. In adapting an ASR task involving human-computer interactions, the proposed approach is found to be as effective as the existing techniques but with a substantial reduction in the computational cost.

© 2016 Elsevier Ltd. All rights reserved.

*Keywords:* Fast adaptation; Acoustic model interpolation; Sparse coding; Exemplar and learned speaker dictionary

## 1. Introduction

The automatic speech recognition (ASR) systems are traditionally developed using the Gaussian mixture model (GMM) based context-dependent hidden Markov model (CD-HMM). Recently, the deep neural network (DNN) (Hinton et al., 2012) has been proposed to overcome the inefficiency of the GMM in the modeling of the events that lie on or near a nonlinear manifold in the data. With the advent of the DNN for generating the observation probabilities, the ASR systems based on DNN-HMM are fast becoming popular<sup>1</sup>. In general, the ASR systems are trained on speech data from a large number of speakers (male and female). This pooled training is primarily intended towards

\* Corresponding author.

E-mail address: [s.syed@iitg.ernet.in](mailto:s.syed@iitg.ernet.in) (S. Shahnawazuddin), [rsinha@iitg.ernet.in](mailto:rsinha@iitg.ernet.in) (R. Sinha).

<sup>1</sup> Even though, the work presented in this paper, while being targeted to the GMM-HMM systems, remains relevant as some of the classical GMM-HMM adaptation approaches have been ported to the DNN-based systems.

8 making the statistical models speaker independent (SI). Unlike a speaker dependent (SD) system, the SI systems  
9 have to deal with both the intra-speaker and the inter-speaker variability. As a result of that, an SI system is reported  
10 to be 2 to 3 times inferior in comparison to an SD system when both the systems are trained with an equal amount of  
11 data (Woodland, 2001). The SD systems, though quite effective, are infeasible to be built for each of the speakers in  
12 the user population on account of requiring a large amount of speech data per speaker. Consequently, the speaker  
13 adaptation techniques have been developed which intend to modify the parameters of the SI system to better suit a  
14 particular speaker given the limited amount of data from that speaker.

15 Some of the applications of ASR are information retrieval, language learning tools, voice-based search and  
16 entertainment (Wang et al., 2008; Eskenazi, 2009; Schalkwyk et al., 2010; Gray et al., 2014). In the case of ASR  
17 tasks involving human-machine interactions, the system is required to recognize speech from the adult as well as the  
18 child (male/female) speakers. Further in such ASR tasks, the adaptation data is made available to the ASR system in  
19 an incremental manner. Moreover, the duration of the data is generally very small. Consequently, the conventional  
20 adaptation techniques like the maximum *a-posteriori* (MAP) (Gauvain and Lee, 1994) and the maximum likelihood  
21 linear regression (MLLR) (Leggetter and Woodland, 1995; Digalakis et al., 1995) are found to be unsuitable. This is  
22 so because the available adaptation data is insufficient to estimate a large number of parameters. On the other hand,  
23 the fast adaptation techniques try to modify the model parameters even with a small amount of adaptation  
24 data (Hazen and Glass, 1997; Gales, 1999; Kuhn et al., 2000; Kenny et al., 2005). In these techniques the adapted  
25 model parameters are derived by a linear interpolation of a set of bases (predefined acoustic models) spanning a low  
26 ( $K$ ) dimensional subspace. Consequently, only a few interpolation weights (or the direction coordinates) need to be  
27 estimated. Such a reduction in complexity makes these approaches amenable for aforementioned adaptation tasks.  
28 The fast adaptation approaches are found to be effective for the adaptation of the Gaussian means as well as the mix-  
29 ture-weights (Duchateau et al., 2008; Hahm et al., 2010) of the acoustic model.

30 In the initial works on fast adaptation (Hazen and Glass, 1997; Gales, 1999; Kuhn et al., 2000; Kenny et al.,  
31 2005), the  $K$  bases to be interpolated were kept fixed. The recent works (Mak et al., 2006; Teng et al., 2009) have  
32 shown that an improved recognition performance can be achieved for the given test data by selecting  $K$  bases from a  
33 predefined set of acoustic models. In those approaches, speaker adapted (SA) models corresponding to each of the  $N$   
34 speakers in the training set are derived first. This is usually done by adapting the mean parameters of the SI model  
35 with the speaker-specific data while keeping the covariance matrices and the mixture-weights unchanged. Given the  
36 test data, a set of  $K$  models (bases) are then selected from the  $N$  acoustic models and are linearly combined. The  $K$   
37 interpolation weights required for linear combination are estimated using an iterative maximum-likelihood (ML)  
38 approach. It is worth mentioning here that there do exist some fast adaptation approaches in which the bases are  
39 derived using the given adaptation data rather than using predefined ones (Gales, 1999; Kenny et al., 2005).

40 The aforementioned fast adaptation techniques differ mainly in the way the  $K$  bases are created/selected. The  
41 work presented in Mak et al. (2006) employs a Vitebi-alignment-based ML search for the basis selection. The ML  
42 search becomes very cumbersome when  $N$  is large and hence is not suitable for adapting interactive ASR tasks. The  
43 approach reported in Teng et al. (2009) requires the estimation of the interpolation weights for all the  $N$  models using  
44 a single iteration of ML estimation. The  $K$  acoustic models having the largest magnitude for the interpolation  
45 weights are selected for deriving the adapted model. The interpolation weights for those  $K$  models are then re-esti-  
46 mated iteratively. The complexity introduced in estimating weights for all the  $N$  models limits the feasibility of this  
47 approach in the case of ASR tasks involving human-machine interactions.

48 The work presented in this paper is inspired by the recent works employing a dynamic selection of acoustically  
49 close bases. In our earlier work (Shahnawazuddin and Sinha, 2014a), techniques employing the sparse representation  
50 (SR) (Elad, 2010) over a redundant dictionary were explored for the dynamic selection of the bases. In that work, the  
51 sparse coding was used only for the basis selection while the basis coefficients in the sparse coding were discarded.  
52 Like the other model-interpolation-based techniques, the interpolation weights were estimated iteratively in the ML  
53 sense. The latency in the basis selection process was reduced due to the greedy SR-based approaches but the compu-  
54 tational cost of the weight estimation remained the same. In contrast to our earlier work, in this paper we propose a  
55 novel use of the sparse coding to derive the adapted model mean parameters.<sup>2</sup> The main contributions of the work  
56 are as follows:

<sup>2</sup> An initial version of this work was presented in *INTERSPEECH 2014* (Shahnawazuddin and Sinha, 2014b)

Download English Version:

<https://daneshyari.com/en/article/4973693>

Download Persian Version:

<https://daneshyari.com/article/4973693>

[Daneshyari.com](https://daneshyari.com)