

Jacobian Approach to Joint Adaptation to Noise, Channel and Vocal Tract Length

Shigeki Sagayama^{1,2}, Yutaka Kato², Mitsuru Nakai² and Hiroshi Shimodaira²

¹ The University of Tokyo, Bunkyo-ku, Tokyo 113-8656 Japan

² Japan Advanced Institute of Science and Technology, Tatsu-no-kuchi, Ishikawa 923-1292 Japan

E-mail: sagayama@hil.t.u-tokyo.ac.jp, {sagayama,ykatou,mit,sim}@jaist.ac.jp

Abstract

This paper describes the Jacobian approach to simultaneously adapting acoustic models to unknown noise, channel and vocal tract length from a supervised adaptation data. As has been both theoretically and experimentally shown, Jacobian adaptation is one of most efficient methods for model adaptation if the target condition is close to the initial condition. It utilizes the linear relationship in the neighbor of the initial condition which in turn can be used in decomposition of multiple factors. The analytic relationship between noise, channel, vocal tract length and the observed cepstrum is linearized using the Jacobian matrices. Least squares fit gives the estimates of noise, channel and vocal tract stretch parameters. Experimental evaluation gave a significant improvement to the recognition accuracy.

1 Introduction

Jacobian approach[4, 5] (first conceived in 1992, experimentally evaluated in 1995, and presented in Japanese in 1996) is a method of adapting initial acoustic models under an initial condition, A , to a target condition, B , assuming that the difference between the two conditions is relatively small. This framework is motivated by speaker adaptation from the initial speaker, A (or speaker-independent), to the target speaker, B . Based on this idea, we introduced an approach that uses a Jacobian matrix. Being a straight forward linear computation, Jacobian adaptation is computationally very efficient and is one of most suitable solutions in practical applications of speech recognition where environmental conditions vary from time to time (e.g., mobile applications) or with each usage (e.g., telephone applications) and cause acoustic model mismatch resulting in a serious degradation of performance. Retraining acoustic models using the actual noise environment is not feasible in most cases in terms of training data amount and computational time constraint. Composing mathematically a noisy speech model from acoustic models of clean speech and noise such as Parallel Model Composition (PMC)[1] and Noise-Voice Composition (NOVO)[2] is often computationally too expensive to follow in real-time the instantaneous changes in noise spectrum and level.

In recent further research results using Jacobian adaptation such as initial model selection[6] and transformation of Jacobian matrices[7, 8], this approach has been made not only more computationally efficient but also more effective in its performance to a given noisy condition.

Differential approach including Jacobian adaptation has a high potential in model adaptation. It can be applied to any factors whose relationship to the feature vectors is analytically known. If the initial and target conditions are close enough for good linear approximation, this approach is applicable to an arbitrary function, not being limited within a specific function such as in the Vector Taylor Series[3] which is related to Jacobian adaptation in a sense that first-order approximation of Taylor series is Jacobian. The Jacobian approach opens up a new non-parametric analytic paradigm of acoustic model adaptation other than existing probabilistic (e.g., Bayesian approaches) or parametric approaches (e.g., MLLR). In addition, by introducing any parametric model, Jacobian approach turns to be a parametric adaptation method.

In this paper, we discuss an extension of the Jacobian approach to simultaneous adaptation to unknown noise and channel, and also to adaptation to vocal tract length.

If there is a small change in linear channel characteristics, this approach is also easily applicable. Given a small change included in the adaptation problem, it is also solved in the same way as in the noise case. It is further extended to adaptation to vocal tract length changes since there is an analytic relationship between vocal tract length stretch and the resulted spectrum.

It is often the case that channel and noise change simultaneously and can not be observed separately. This arouses a joint estimation problem of channel and noise. This may be solved with a computationally expensive iterative algorithm. If, however, the changes of channel and noise characteristics between the initial and target conditions are small, Jacobian approach can make the problem linear so that simple least squares fit gives a computationally efficient algorithm.

In this paper, we extend the formulation to include the channel problem. We also discuss the future directions of this approach.

2 Jacobian Adaptation to Noise and Channel

2.1 Formulating Jacobian Adaptation

If a vector variable, C_Y , is an analytic function of other variables, C_S , C_N and C_H , namely,

$$C_Y = \Psi(C_S, C_N, C_H), \quad (1)$$

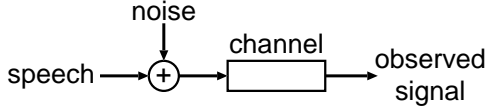


Figure 1: The model of speech, noise and channel.

a small change in C_Y caused by small changes in C_S , C_N and C_H is given by

$$\Delta C_Y = \frac{\partial \Psi}{\partial C_S} \Delta C_S + \frac{\partial \Psi}{\partial C_N} \Delta C_N + \frac{\partial \Psi}{\partial C_H} \Delta C_H \quad (2)$$

regardless of the meaning of variables C_S , C_N and C_H . We call $\frac{\partial \Psi}{\partial C_S}$, $\frac{\partial \Psi}{\partial C_N}$ and $\frac{\partial \Psi}{\partial C_H}$ Jacobian matrices whose (i, j) component is the derivative of the i th component of Ψ in respect to the j th component of C , i.e., $\frac{\partial \Psi_i}{\partial C_j}$.

The above mathematical relationship holds if n -dimensional vectors S , N , H , and Y represents clean speech spectrum, additive noise spectrum, multiplicative channel characteristics (transfer function in the power spectral domain), and the resulted composite speech spectrum, respectively, and C_S , C_N , C_H and C_Y represent their corresponding cepstrum vectors.

When noise and channel conditions represented by C_N and C_H ("Condition A") change into $C_N + \Delta C_N$ and $C_H + \Delta C_H$ ("Condition B") with the speech spectrum fixed, the composite cepstrum also changes into:

$$C_Y + \Delta C_Y = \Psi(C_S, C_N + \Delta C_N, C_H + \Delta C_H) \quad (3)$$

where ΔC_H is given by Eq.(2) and $\Delta C_S = 0$. This is the basic idea of Jacobian adaptation to a new condition.

2.2 Jacobian for Noise and Channel

The relationship among n -dimensional vectors, S , N , H , and Y shown in Fig. 1 is given by

$$Y = H(S + N) \quad (4)$$

in the linear spectral domain. These relations also hold in a particular case where C_S , C_N , C_H and C_Y represent cepstra of speech, noise, channel and resulted composite signal, respectively, as shown in In the linear spectral domain, holds Equivalently, in the cepstral domain, we have¹

$$C_Y = F^* \left[\log \{ \exp(F C_S) + \exp(F C_N) \} \right] + C_H \quad (5)$$

where F is the Fourier transform matrix² and F^* is the transposed complex conjugate of F that $F^* F = \mathbf{1}$.

The signal-to-noise ratio is not treated separately; the noise power is included in N . Since these cepstra are related to the respective spectra by

$$\begin{aligned} \log Y = F C_Y, \quad \log S = F C_S, \\ \log N = F C_N, \quad \log H = F C_H \end{aligned} \quad (6)$$

¹For vectors \mathbf{a} and \mathbf{b} , we define $\mathbf{a}\mathbf{b} = (a_1 b_1, \dots, a_n b_n)^T$, $\mathbf{a}/\mathbf{b} = (a_1/b_1, \dots, a_n/b_n)^T$, $\exp \mathbf{a} = (\exp a_1, \dots, \exp a_n)^T$, and $\log \mathbf{a} = (\log a_1, \dots, \log a_n)^T$. One can regard vectors as diagonal matrices for consistency with matrix arithmetic.

²Fourier transform matrix, F , is not necessarily a square matrix. It depends on the resolution of frequency and the number of cepstral points (qufrecencies).

It should be noted that the Fourier transform is essentially the cosine transform in real symmetric spectrum cases, namely, $F_{ij} = \cos \frac{2ij\pi}{n}$ and the computation amount can be reduced to one quarter by handling the positive frequency only. One implementation is

$$F_{ik} = \cos \frac{i(k+0.5)\pi}{N} \quad (7)$$

where N stands for the number of frequency points. Cepstrum can be simply replaced by MFCC in the above formulation where the corresponding spectrum is replaced by mel-frequency warped spectrum.

If these changes are small, the resulted change, ΔC_Y , is

$$\Delta C_Y = \frac{\partial C_Y}{\partial C_N} \Delta C_N + \Delta C_H \quad (8)$$

according to Eq. (2).

The Jacobian matrix is easily calculated at the initial condition A:

$$\begin{aligned} J_C &\equiv \frac{\partial C_Y}{\partial C_N} = \frac{\partial C_Y}{\partial \log Y} \frac{\partial \log Y}{\partial Y} \frac{\partial Y}{\partial N} \frac{\partial N}{\partial \log N} \frac{\partial \log N}{\partial C_N} \\ &= F^* \frac{\mathbf{1}}{H(S+N)} H N F = F^* \frac{N}{S+N} F \end{aligned} \quad (9)$$

Thus, if differences between the initial and observed conditions, A and B, is found in the cepstrum domain, i.e., ΔC_N and ΔC_H , the composite cepstrum, $C_Y + \Delta C_Y$, is approximately computed by Eq. (8).

Eq.(9) implies that any analytic relationship can be decomposed into a product and/or sum of simple Jacobian matrices. Just like Bayesian networks, we can consider "Jacobian newtworks" which include multiple causes and intermediate results connected in a network structure.

2.3 Least Mean Squares Estimation

Noise and channel differences, however, can not be observed separately in most cases such as in telephone speech recognition. It is therefore necessary to simultaneously estimate the noise and channel spectral characteristics from the observed speech signal.

Suppose that we observe noisy speech containing multiple phonemes through a channel. Note that Eq. (4) commonly holds at all different S corresponding to multiple phonemes. If M hidden states of the model are aligned (e.g., by Viterbi algorithm) to the input speech, we have a set of M different simultaneous relations derived from Eq. (4). Assuming that the target condition B is relatively close to the initial condition A, we have a set of M linear equations with additional errors derived from Eq. (8) as follows:

$$\begin{cases} \Delta C_Y^{(1)} = J_C^{(1)} \Delta C_N + \Delta C_H + \epsilon^{(1)} \\ \Delta C_Y^{(2)} = J_C^{(2)} \Delta C_N + \Delta C_H + \epsilon^{(2)} \\ \vdots \\ \Delta C_Y^{(M)} = J_C^{(M)} \Delta C_N + \Delta C_H + \epsilon^{(M)} \end{cases} \quad (10)$$

Thus, if a small amount of input speech with unknown noise and channel is given with its phonetical transcription, we can obtain the joint estimate of ΔC_N and ΔC_H by simple

Table 1: Experimental conditions for noise and channel simultaneous adaptation.

Speech DB	ATR Speech DB A-set (5240 words)
Speakers	MAU (male), FFS (female)
Training	2650 words (odd-numbered)
Testing	655 words (from even-numbered)
Features	13 MFCCs + 13 Δ MFCCs
Models	3-state, 4-mixture, CD phone HMMs
Noise A	Exhibition Hall (10dB SNR)
Noise B	Crowd (0, 10, 20, 30 dB)
Channel A	Flat
Channel B	Simulated (shown right)

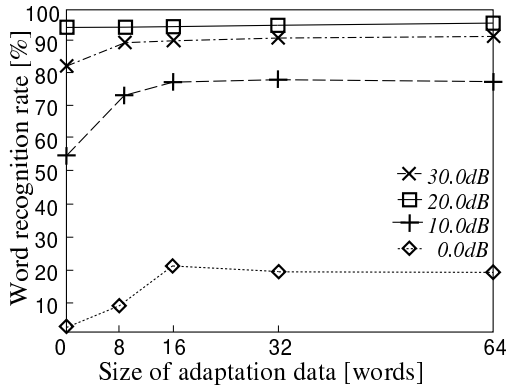


Figure 2: A typical example of simultaneous noise and channel adaptation of 10-dB SNR initial models to given noise and channel conditions.

least squares estimation minimizing the sum of squares of error terms, $\sum |\epsilon^{(i)}|^2$. This procedure can be iterated for more accurate Viterbi alignment using the adapted models.

This process estimates only two vectors, ΔC_N and ΔC_H , and is advantageous from the statistical estimation point of view. Once the linear decomposition into ΔC_N and ΔC_H are estimated, Eq. (8) is applied to all mean vectors of the model for adapting the whole model.

2.4 Experimental Evaluation

For experimental evaluation of simultaneous adaptation to noise and channel, speaker-dependent isolated-word recognition was performed with artificially added environmental noise and simulated channel under conditions shown in Table 1. A typical result is shown in Fig. 2. Table 2 shows the averaged results from adaptation experiments across different noise sources and SNRs ranging from 0dB to 30dB. A significant error reduction around 10% is observed for only 8 words given for supervised adaptation, while more data do not give more improvement. This is due to the quite limited number of free parameters to estimate in the least squares fit in the linear decomposition of mixed factors.

Table 2: Simple average recognition and error reduction rates in simultaneous adaptation across mismatched noise conditions.

Adapt. Data Size (wds)	0	8	16	32	64
Ave. Recog. Rate (%)	33.6	40.6	40.9	40.3	40.0
Error Reduc. Rate (%)	—	10.6	10.6	10.1	9.8

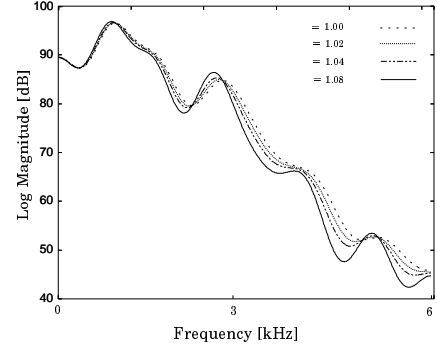


Figure 3: Resulted spectrum of the cepstral-domain frequency stretching by Jacobian approximation

3 Jacobian Adaptation to Vocal Tract Length

3.1 Vocal Tract λ -stretched Cepstrum

The Jacobian approach can be extended to include some aspects of speaker differences. If the vocal tract length becomes λ -times longer, the corresponding speech spectrum changes from $f(\omega)$ into $f(\lambda\omega)$. We can again calculate the Jacobian matrix J_λ of resulted cepstrum in respect to λ . (This matrix has only one column and appears like a vector.)

Combining noise, channel, and vocal tract length factors, we can express the small changes in the composite cepstrum as follows using small changes in noise cepstrum, channel cepstrum, and the vocal length stretch coefficient:

$$\Delta C_Y = J_N \Delta C_N + \Delta C_H + J_\lambda \Delta \lambda \quad (11)$$

In Eqs.(6), if we replace the Fourier Transform F with a λ -stretched Fourier transform F^λ , or, if we use

$$F_{ik}^\lambda = \cos \frac{\lambda i(k + 0.5)\pi}{N} \quad (12)$$

instead of Eq.(7), the λ -stretched speech spectrum \tilde{S} is given by

$$\log \tilde{S} = F^\lambda C_S. \quad (13)$$

The λ -stretched cepstrum \tilde{C}_S is thus expressed as

$$\tilde{C}_S = F^{-1} F^\lambda C_S \quad (14)$$

whose i -th component is given by

$$\tilde{C}_{Si} = \sum_{j=1}^N F_{ij}^{-1} \sum_{k=1}^p F_{jk}^\lambda C_{Sk} \quad (15)$$

from which the i -th component of its Jacobian in respect to λ is derived by differentiating it by λ as

$$J_{\lambda i} = \sum_{j=1}^N F_{ij}^{-1} \sum_{k=1}^p \frac{-j(k + 0.5)\pi}{N} G_{jk} C_{Sk} \quad (16)$$

where matrix G represents the sine transform.

Fig. 3 shows the spectral domain display of cepstral-domain frequency stretching with Jacobian linear approximation by $J_\lambda \Delta\lambda$. From this figure, we can understand that the approximation works only near $\lambda = 1.0$ but has a limitation as λ leaves 1.0.

3.2 Joint Adaptation to Three Factors

As preliminary evaluation of frequency stretching, we calculated appropriate λ values by comparing spectral peak frequencies of initial model (condition A) and target (condition B) spectra derived from observed cepstra. Among several ways of finding and applying λ values to the adaptation procedure in a speaker adaptation experiments across 4 speakers, averaged λ values extracted from vowels and applied to all phonemes yielded 4% recognition error reduction with 2 training words and 5% with 4, 8, and 16 words.

In joint adaptation to noise, channel and speaker, the following procedure was taken:

- Preparation (before adaptation):
 1. Train the initial noise model (C_N^A)
 2. Train the noisy speech model (C_Y^A)
 3. Calculate Jacobian matrices (J_N, J_λ)
- Adaptation (supervised):
 4. Train the target noisy speech model (C_Y^B) using a small amount of supervised training data through Viterbi time alignment.
 5. Observe the difference (ΔC_Y^B) between conditions A and B .
 6. LMS-estimate the differences of noise and channel ($\Delta C_N, \Delta C_H$) between A and B .
 7. Find $\Delta\lambda$ by comparing C_Y^A and C_Y^B .
 8. Calculate $\Delta \tilde{C}_Y^A$ using $\Delta\lambda$ and J_λ .
 9. Obtain the adapted model (\tilde{C}_Y^B) combining the changes $\Delta C_N, \Delta C_H$ and ΔC_λ .
- Recognition:
 10. Recognize the speech data in condition B using the adapted model.

A typical results of applying this joint adaptation technique to different noise, channel and speaker is shown in Fig. 4 where 16 LPC cepstrum coefficients and their time derivatives were used as features. It was experimentally found that joint adaptation gives more effective results than the sum of separate adaptation to separate factors. A possible explanation is that speaker differences include not only vocal tract length stretch but also multiplicative (and possibly additive) spectrum differences, and that the mixed problem of noise, channel and speaker adaptation is better solved by joint adaptation in respect to all of three factors.

4 Conclusion

This paper introduced linear decomposition of noise and channel differences in mismatched conditions using a Jacobian formulation. Joint adaptation to noise, channel and

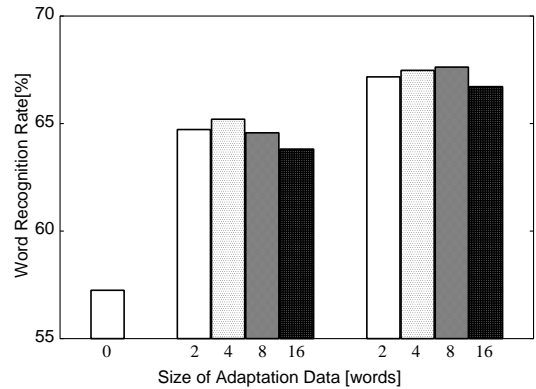


Figure 4: A typical results of word speech recognition of no adaptation (left), joint adaptation to additive (noise), multiplicative (channel) factors (center) and adaptation to additive, multiplicative and frequency stretching (vocal tract length) factors. Initial condition (A): models were trained with 3 speakers (mht, mms, mmy), flat channel, 30dB SNR noise “in factory”; Testing condition (B): speaker “mau”, low-pass channel, 10dB SNR noise “crossroads”.

vocal tract length (frequency stretching) mismatches and demonstrated with experimental results was also discussed. Future works include overall least mean squares joint estimation of noise, channel, and vocal tract length formulated in a similar way as described for the noise and channel case.

References

- [1] M. J. F. Gales and S. J. Young, “An Improved Approach to the Hidden Markov Model Decomposition of Speech and Noise,” *Proc. ICASSP92*, pp. 233–236, 1992.
- [2] F. Martin, et al., “Recognition of Noisy Speech by Using the Composition of Hidden Markov Models,” *Proc. 1992 Fall ASJ Conf.*, 1-7-10, 1992.
- [3] P. J. Moreno, B. Raj, and R. Stern, “A Vector Taylor Series Approach for Environment-Independent Speech Recognition,” *Proc. ICASSP96*, pp. 733–736, 1996.
- [4] S. Sagayama, Y. Yamaguchi, S. Takahashi, and J. Takahashi: “Jacobian Approach to Fast Acoustic Model Adaptation,” *Proc. ICASSP97*, Vol. 2, pp. 835–838, 1997.
- [5] Y. Yamaguchi, S. Takahashi, and S. Sagayama, “Fast Adaptation of Acoustic Models to Environmental Noise Using Jacobian Adaptation Algorithm,” *Proc. Eurospeech 97*, pp. 2051–2054, 1997.
- [6] H. Shimodaira, T. Akae, M. Nakai and S. Sagayama, “Jacobian Adaptation of HMM with Initial Model Selection for Noisy Speech Recognition,” *Proc. ICSLP2000 (Beijing)*, pp.1003-1006, 2000.
- [7] C. Cerisara, L. Rigazio, R. Boman and J.-C. Junqua, “Transformation of Jacobian Matrices for Noisy Speech Recognition,” *Proc. ICSLP2000 (Beijing)*, pp. 369-372, 2000.
- [8] R. Sarikaya and J. Hansen, “Improved Jacobian Adaptation for Fast Acoustic Model Adaptation in Noisy Speech Recognition,” *Proc. ICSLP2000 (Beijing)*, pp. 702–705, 2000.