



# Analysis of inter-articulator correlation in acoustic-to-articulatory inversion using generalized smoothness criterion

Prasanta Kumar Ghosh and Shrikanth S. Narayanan

Signal Analysis and Interpretation Laboratory (SAIL), Department of Electrical Engineering,  
University of Southern California, Los Angeles, CA 90089

prasantg@usc.edu, shri@sipi.usc.edu

## Abstract

The movements of the different speech articulators are known to be correlated to various degrees during speech production. In this paper, we investigate whether the inter-articulator correlation is preserved among the articulators estimated through acoustic-to-articulatory inversion using the generalized smoothness criterion (GSC). GSC estimates each articulator separately without explicitly using any correlation information between the articulators. Theoretical analysis of inter-articulator correlation in GSC reveals that the correlation between any two estimated articulators may not be identical to that between the corresponding measured articulatory trajectories; however, based on smoothness constraints provided by the real articulatory data, we found that, in practice, the correlation among articulators is approximately preserved in GSC based inversion. To validate the theoretical analysis on inter-articulator correlation, we propose a modified version of GSC where correlations among articulators are explicitly imposed. We found that there is no significant benefit in inversion using such modified GSC, which further strengthens the conclusions drawn from the theoretical analysis of inter-articulator correlation.

**Index Terms:** acoustic-to-articulatory inversion, inter-articulation correlation, generalized smoothness criterion

## 1. Introduction

Estimation of representations in the articulatory space from representations in the acoustic space is known as acoustic-to-articulatory inversion. In this paper, we consider the problem of estimating articulatory position vector sequence (or trajectory) from a given MFCC (Mel Frequency Cepstral Coefficients of short time acoustic speech signal frames) vector sequence. Exemplary articulatory vectors can correspond to the positions (i.e., X and Y values in mid-sagittal plane) of the upper lip (UL), lower lip (LL), lower incisor (LI), tongue tip (TT), tongue body (TB), tongue dorsum (TD), and the velum (V), as considered in this paper.

Previous studies have provided evidence that the mapping between the acoustic and articulatory spaces is not unique [1]. Therefore, many inversion techniques [2, 3, 4, 5, 6, 7, 8, 9] have been proposed where issues of non-uniqueness are mitigated by constraining the articulator trajectories to be smooth – e.g., low-pass filtering the estimated articulator trajectories [7, 8, 9]. We recently introduced a technique for incorporating the smoothness constraint inside the optimization problem using generalized smoothness criterion (GSC) [2]. In contrast to other inversion techniques, GSC estimates each articulator trajectory in an independent fashion using the acoustic feature sequence from a test utterance. However, it is well known that many of the measured articulators' movements can be correlated with one another [10]. For example, the upper and lower lips' movements are correlated

since they move together to create lip opening and closure. Similarly since TT, TB, TD are three locations on the same physical tongue organ and, hence, their movements are expected to be correlated too. It is however not clear whether the GSC formulation preserves the inter-articulator correlation since it estimates each articulator independently.

In this work, we perform a theoretical analysis on the correlations among articulators estimated using GSC and derive their relations to the correlations among measured articulatory trajectories. Based on the analysis of inter-articulator correlation in GSC we show that, theoretically, there is no guarantee that GSC preserves the inter-articulator correlation as observed in the training data. However, when the theoretical analysis is examined with respect to real articulatory data, we found that the differences between the correlations among estimated articulators and those among measured articulators are not significant. Thus, it turns out that, in practice, the correlation among articulators estimated by GSC are approximately similar to those among measured articulators.

To further validate this inter-articulator correlation property of GSC, we develop a way within the GSC framework by which inter-articulator correlation can be explicitly imposed in the acoustic-to-articulatory inversion such that the empirical correlation coefficient between any two estimated articulatory trajectories will be identical to that between the respective articulatory variables (measured) in training data. The accuracy of the inversion obtained by exploiting the inter-articulator correlation is compared experimentally against the accuracy obtained by treating each articulator independently during inversion using GSC [2]. Based on the comparison, we observe that there is no significant benefit in explicitly imposing correlation in the inversion using GSC, which further justifies the validity of the theoretical analysis.

## 2. Dataset and pre-processing

For the analysis and experiments of this paper, we use the Multi-channel Articulatory (MOCHA) database [11] containing acoustic and corresponding articulatory ElectroMagnetic Articulography (EMA) data from one male and one female subject. The articulatory position data have high frequency noise resulting from EMA measurement error. Also the mean position of the articulators changes from utterance to utterance; hence, the position data needs pre-processing before it can be used for analysis. Following the pre-processing steps as outlined by Ghosh et al. [2], we obtain parallel acoustic and articulatory data at a frame rate of 100 observations per second. Of the 460 utterances available from each subject, data from 368 utterances (80%) are used for training, 37 utterances (8%) as the development set (dev set), and the remaining 55 utterances (12%) as the test set.

<sup>†</sup>Work supported by NIH and NSF.

### 3. Generalized smoothness criterion (GSC) for articulatory inversion

In this section, we briefly describe the principle of GSC [2] for inversion. Let  $\{(\mathbf{z}_i, \mathbf{x}_i); 1 \leq i \leq T\}$  represent the parallel acoustic feature vector and articulatory position vector pairs in the training set, where  $\mathbf{z}_i$  and  $\mathbf{x}_i$  represent the  $i^{\text{th}}$  acoustic and articulatory vector respectively.  $\mathbf{x}_n = [x_n^1 \ x_n^2 \ \dots \ x_n^{14}]^T$ , where the 14 elements correspond to X and Y co-ordinates of seven articulators considered in this work.  $[\cdot]^T$  denotes the transpose operator. Now suppose, for the acoustic-to-articulatory inversion, a (test) speech utterance is given and the acoustic feature vectors computed for this utterance are denoted by  $\{\mathbf{u}_n; 1 \leq n \leq N\}$ . The GSC is used to estimate the  $j^{\text{th}}$  articulatory position trajectory  $\{x_n^j; 1 \leq n \leq N\}$  by solving the following optimization problem [2]:

$$\arg \min_{\{x_n^j\}} \left[ \sum_n \left( y^j[n] \right)^2 + C_j \sum_n \sum_l \left( x_n^j - \eta_n^{l,j} \right)^2 p_n^l \right], \quad (1)$$

where,  $y_n^j = \sum_{k=1}^N x_k^j h_{n-k}^j$  and  $h_n^j$  is an articulator specific high-pass filter with cut-off frequency  $f_c^j$ . Thus the first term on the right hand side of Eq. (1) is used to minimize the high frequency components in  $x_n^j$  so that the articulatory position trajectory is smooth.  $\{\eta_n^{l,j}; 1 \leq l \leq L\}$  are the  $L$  possible values of the  $j^{\text{th}}$  articulatory position at the  $n^{\text{th}}$  frame of the test speech utterance and  $p_n^l$  are their probabilities [2].  $C_j$  is the trade off parameter between two terms in the objective function (Eq. (1)).

The solution of Eq. (1) for the  $j^{\text{th}}$  articulator can be written as

$$\mathbf{x}^{j*} = \left( \mathbf{R}^j + C_j \mathbf{I} \right)^{-1} \mathbf{d}^j, \quad (2)$$

where,  $\mathbf{x}^{j*} = [x_1^{j*} \ \dots \ x_N^{j*}]^T$  is the optimal articulatory trajectory.  $\mathbf{R}^j = \{R_{pq}^j\} = \{R^j(p-q)\} \triangleq \sum_n h^j[n-p]h^j[n-q]$ .  $\mathbf{I}$  is the identity matrix and  $\mathbf{d}^j = [C_j \sum_l \eta_n^{l,j} p_n^l, \dots, C_j \sum_l \eta_n^{l,j} p_n^l]^T$ . Note that the solution (Eq. (2)) can be obtained recursively with frame index  $n$  without any loss in accuracy [2].

### 4. Correlation among estimated articulator trajectories

We theoretically investigate how the correlations among estimated articulatory trajectories in GSC differ from those among the measured articulatory trajectories.

From Eq. (2), we can write (for  $i \neq j$ ) the estimated articulatory trajectory in the following way:

$$\begin{cases} x_n^{j*} = \sum_{m_1=1}^N a_{m_1} C_j \sum_{l=1}^L \eta_{m_1}^{l,j} p_{m_1}^l \\ x_n^{i*} = \sum_{m_2=1}^N b_{m_2} C_i \sum_{l=1}^L \eta_{m_2}^{l,i} p_{m_2}^l \end{cases} \quad (3)$$

where,  $\{a_{m_1}, 1 \leq m_1 \leq N\}$  and  $\{b_{m_2}, 1 \leq m_2 \leq N\}$  are the  $n^{\text{th}}$  rows of  $(\mathbf{R}^j + C_j \mathbf{I})^{-1}$  and  $(\mathbf{R}^i + C_i \mathbf{I})^{-1}$  respectively. Let us assume  $\eta_n^l$ , i.e.,  $[\eta_n^{1,1}, \dots, \eta_n^{1,14}]$ , is *i.i.d.* random vector  $\forall l, n$ . Also let the  $j^{\text{th}}$  and  $i^{\text{th}}$  articulators, i.e.,  $\eta_n^{l,j}$  &  $\eta_n^{l,i}$ , have correlation co-efficients  $\rho_{ji}$ . In other words,

$$\begin{aligned} \mathbb{E} \left[ \eta_n^{l,j} \right] &= \mu_j, & \mathbb{V} \left[ \eta_n^{l,j} \right] &= \mathbb{E} \left[ \left( \eta_n^{l,j} - \mu_j \right)^2 \right] = \sigma_j^2, \\ \text{COV} \left( \eta_{n_1}^{l_1,j}, \eta_{n_2}^{l_2,i} \right) &= \mathbb{E} \left[ \left( \eta_{n_1}^{l_1,j} - \mu_j \right) \left( \eta_{n_2}^{l_2,i} - \mu_i \right) \right] \\ &= \rho_{ji} \sigma_i \sigma_j \delta_{n_1, n_2} \delta_{l_1, l_2}, \end{aligned} \quad (4)$$

where,  $\mathbb{E}[\cdot]$ ,  $\mathbb{V}[\cdot]$ , and  $\text{COV}[\cdot, \cdot]$  denote the mean, variance, and covariance of the random variables.  $\delta_{m,n}$  is the kronecker delta, i.e.,  $\delta_{m,n} = 1$ , when  $m = n$  and  $\delta_{m,n} = 0$ , when  $m \neq n$ . Therefore, the mean of  $x_n^{j*}$  and  $x_n^{i*}$  are as follows:

$$\begin{aligned} \mathbb{E} \left[ x_n^{j*} \right] &= \sum_{m_1=1}^N a_{m_1} C_j \sum_{l=1}^L \mathbb{E} \left[ \eta_{m_1}^{l,j} \right] p_{m_1}^l \\ &= \mu_j C_j \sum_{m_1=1}^N a_{m_1} \sum_{l=1}^L p_{m_1}^l \\ &= \mu_j C_j \sum_{m_1=1}^N a_{m_1} \left( \sum_{l=1}^L p_{m_1}^l = 1 \right) \end{aligned} \quad (5)$$

$$\text{Similarly, } \mathbb{E} \left[ x_n^{i*} \right] = \mu_i C_i \sum_{m_2=1}^N b_{m_2} \quad (6)$$

Therefore,

$$\begin{aligned} \mathbb{V} \left[ x_n^{j*} \right] &= \mathbb{E} \left[ \left( x_n^{j*} - \mathbb{E} \left[ x_n^{j*} \right] \right)^2 \right] \\ &= C_j^2 \left( \sum_{m_{11}=1}^N \sum_{m_{12}=1}^N \sum_{l_1=1}^L \sum_{l_2=1}^L a_{m_{11}} a_{m_{12}} p_{m_{11}}^{l_1} p_{m_{12}}^{l_2} \right. \\ &\quad \left. \mathbb{E} \left[ \left( \eta_{m_{11}}^{l_1,j} - \mu_j \right) \left( \eta_{m_{12}}^{l_2,j} - \mu_j \right) \right] \right) \\ &= C_j^2 \left( \sum_{m_{11}=1}^N \sum_{m_{12}=1}^N \sum_{l_1=1}^L \sum_{l_2=1}^L a_{m_{11}} a_{m_{12}} p_{m_{11}}^{l_1} p_{m_{12}}^{l_2} \right. \\ &\quad \left. \sigma_j^2 \delta_{m_{11}, m_{12}} \delta_{l_1, l_2} \right) \text{ (using Eq. (4))} \\ &= C_j^2 \sigma_j^2 \sum_{m_1=1}^N \sum_{l=1}^L \left( a_{m_1} p_{m_1}^l \right)^2 \end{aligned} \quad (7)$$

$$\& \mathbb{V} \left[ x_n^{i*} \right] = C_i^2 \sigma_i^2 \sum_{m_2=1}^N \sum_{l=1}^L \left( b_{m_2} p_{m_2}^l \right)^2 \quad (8)$$

And, the covariance between the  $j^{\text{th}}$  and  $i^{\text{th}}$  estimated articulators is

$$\begin{aligned} \text{COV} \left( x_n^{j*}, x_n^{i*} \right) &= \mathbb{E} \left[ \left( x_n^{j*} - \mathbb{E} \left[ x_n^{j*} \right] \right) \left( x_n^{i*} - \mathbb{E} \left[ x_n^{i*} \right] \right) \right] \\ &= C_i C_j \left( \sum_{m_1=1}^N \sum_{m_2=1}^N \sum_{l_1=1}^L \sum_{l_2=1}^L a_{m_1} b_{m_2} p_{m_1}^{l_1} p_{m_2}^{l_2} \right. \\ &\quad \left. \mathbb{E} \left[ \left( \eta_{m_1}^{l_1,j} - \mu_j \right) \left( \eta_{m_2}^{l_2,i} - \mu_i \right) \right] \right) \\ &= C_i C_j \left( \sum_{m_1=1}^N \sum_{m_2=1}^N \sum_{l_1=1}^L \sum_{l_2=1}^L a_{m_1} b_{m_2} p_{m_1}^{l_1} p_{m_2}^{l_2} \rho_{ji} \right. \\ &\quad \left. \sigma_i \sigma_j \delta_{m_1, m_2} \delta_{l_1, l_2} \right) \\ &= C_i C_j \rho_{ji} \sigma_i \sigma_j \left( \sum_{m=1}^N \sum_{l=1}^L a_m b_m \left( p_m^l \right)^2 \right) \end{aligned}$$

Hence, the correlation coefficient between the estimated  $j^{\text{th}}$  and  $i^{\text{th}}$  articulators, i.e.,  $x_n^{j*}$  &  $x_n^{i*}$  is

$$\begin{aligned} \rho_{ji}^* &= \frac{\text{COV} \left( x_n^{j*}, x_n^{i*} \right)}{\sqrt{\mathbb{V} \left[ x_n^{j*} \right]} \sqrt{\mathbb{V} \left[ x_n^{i*} \right]}} \\ &= \frac{\sum_{m=1}^N \sum_{l=1}^L a_m b_m \left( p_m^l \right)^2}{\sqrt{\sum_{m=1}^N \sum_{l=1}^L \left( a_m p_m^l \right)^2} \sqrt{\sum_{m=1}^N \sum_{l=1}^L \left( b_m p_m^l \right)^2}} \rho_{ji} \end{aligned} \quad (9)$$

Thus, the correlation  $\rho_{ji}^*$  between the  $j^{\text{th}}$  and  $i^{\text{th}}$  articulator trajectories estimated using GSC is not necessarily identical to the inter-articulator correlation  $\rho_{ji}$  in the training set. It is important to note that when the cut-off frequencies of the  $j^{\text{th}}$  and  $i^{\text{th}}$  articulator-specific filters are nearly identical, their impulse responses (i.e.,  $h_n^j$  and  $h_n^i$ ) and hence the correlation matrices  $\mathbf{R}^j$  and  $\mathbf{R}^i$  are approximately same. In addition, if the trade-off parameters of the respective articulators, i.e.,  $C_j$  and  $C_i$ , are similar, then  $a_m \approx b_m$ ,  $1 \leq m \leq N$  because  $\{a_{m_1}, 1 \leq m_1 \leq N\}$  and  $\{b_{m_2}, 1 \leq m_2 \leq N\}$  are the  $n^{\text{th}}$  rows of  $(\mathbf{R}^j + C_j \mathbf{I})^{-1}$  and  $(\mathbf{R}^i + C_i \mathbf{I})^{-1}$  respectively. Under such circumstances (i.e.,  $a_m \approx b_m$ ), it is easy to see that  $\rho_{ji}^* \approx \rho_{ji}$ ; this means that when  $a_m \approx b_m$ , the correlation is approximately preserved among the articulators estimated by GSC. However, as shown by Ghosh et al. [2], neither  $f_c^j$  nor  $C_j$  is identical across different articulators and, hence, we compute  $\rho_{ji}^*/\rho_{ji}$  (from Eq. (9)) for each test utterance. Fig. 1 demonstrates the average values of  $\rho_{ji}^*/\rho_{ji}$  for each pair of articulatory variables ( $j$ -th and  $i$ -th,  $1 \leq j, i \leq 14$ ) separately for each subject in the MOCHA database. Note that standard deviation (SD) of  $\rho_{ji}^*/\rho_{ji}$  over all test utterances is minimal (maximum SD among all pair of articulatory variables is  $5.54 \times 10^{-5}$ ).

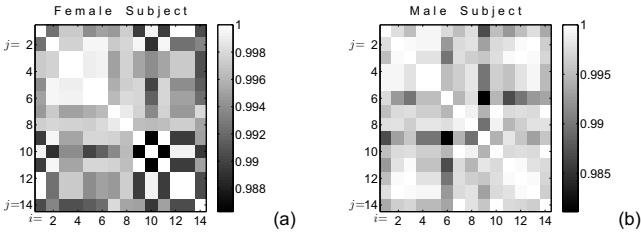


Figure 1:  $\rho_{ji}^*/\rho_{ji}$  for all pairs of articulators for (a) female (b) male subjects in the MOCHA database.  $i$  and  $j$  vary over articulatory variable index  $1, \dots, 14$ .

It is clear from Fig. 1 that  $\rho_{ji}^*/\rho_{ji}$  is very close to 1 for all pairs of articulators and hence  $\rho_{ji}^*$  is approximately same as  $\rho_{ji}$  for each subject in the MOCHA database. Therefore, in practice GSC approximately preserves the inter-articulator correlation although theoretically the correlation among estimated articulators is not identical to that among measured articulatory variables.

To further validate this conclusion based on this theoretical analysis, we develop a modified version of GSC framework where inter-articulator correlation is explicitly imposed among estimated articulators during inversion. The goal is to examine where there is any significant benefit in inversion by explicitly preserving the inter-articulator correlation and thereby examine the validity of the theoretical analysis above.

## 5. Modified GSC to preserve inter-articulator correlation

From Eq. (9) it is worth noting that if the variables in the training set ( $\eta_n^{l,j}$  &  $\eta_n^{l,i}$ ) were uncorrelated (i.e.,  $\rho_{ji} = 0$ ), then the estimated trajectories ( $x_n^{j*}$  and  $x_n^{i*}$ ) would also be uncorrelated (i.e.,  $\rho_{ji}^* = 0$ ). This observation motivates us to transform the articulatory position variables,  $\{x_n^j; j = 1, \dots, 14\}$ , into another set of variables,  $\{\tilde{x}_n^j; j = 1, \dots, 14\}$ , where  $\tilde{x}_n^j$  and  $\tilde{x}_n^k$  are uncorrelated  $\forall j, k$ . The GSC can be used in the transformed variable domain for inversion and, after inversion, the correlation between variables can be imposed by transforming them back to the original articulatory position variable domain. Whitening is one of the approaches where a random vector ( $\mathbf{x}_n$ , in our case) is linearly transformed to make its components uncorrelated [12]. We transform  $\mathbf{x}_n$  to obtain  $\tilde{\mathbf{x}}_n$ .

### 5.1. Transformation of the articulatory position vector

Let  $\mu_{\mathbf{x}}$  and  $\mathbf{K}_{\mathbf{xx}}$ , respectively, be the mean vector and the covariance matrix of the random vector  $\mathbf{x}_n$ . Let the eigen decomposition of  $\mathbf{K}_{\mathbf{xx}} = \mathbf{V} \Lambda \mathbf{V}^T$ , where  $\mathbf{V}$  is the orthogonal eigenvector matrix ( $\mathbf{V}^T \mathbf{V} = \mathbf{I}$ ) and  $\Lambda$  is the diagonal eigen value matrix. The following linear transformation whitens  $\mathbf{x}_n$  to the random vector  $\tilde{\mathbf{x}}_n$ ; the components of  $\tilde{\mathbf{x}}_n$  are uncorrelated as it is easy to show that  $\tilde{\mathbf{x}}_n$  has a diagonal covariance matrix  $\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}}$ .

$$\tilde{\mathbf{x}}_n = \mathbf{V}^T \mathbf{x}_n \quad (10)$$

where,  $\mu_{\tilde{\mathbf{x}}} = \mathbf{E}(\tilde{\mathbf{x}}_n) = \mathbf{V}^T \mu_{\mathbf{x}}$ ,  $\mathbf{K}_{\tilde{\mathbf{x}}\tilde{\mathbf{x}}} = \mathbf{V}^T \mathbf{K}_{\mathbf{xx}} \mathbf{V} = \Lambda$ .

Note that the component of  $\tilde{\mathbf{x}}_n$  does not correspond to any physically meaningful articulatory parameters any more. However,  $\mathbf{x}_n$  can be recovered from  $\tilde{\mathbf{x}}_n$  by  $\mathbf{x}_n = (\mathbf{V}^T)^{-1} \tilde{\mathbf{x}}_n = \mathbf{V} \tilde{\mathbf{x}}_n$ .

### 5.2. Frequency analysis of the transformed variables

In the GSC formulation [2], the cut-off frequencies of the articulator specific high-pass filters  $h_n^j$  are determined based on the analysis of the frequency content of all articulatory position variables  $\{x_n^j; j = 1, \dots, 14\}$ . Similarly, to estimate  $\tilde{x}_n^j$  using GSC, we need to analyze the frequency content of the transformed variables  $\{\tilde{x}_n^j; j = 1, \dots, 14\}$  to determine the cut-off frequency of the high-pass filters in the transformed variable domain. We calculate the frequency  $f_c^j$  below which  $\alpha\%$  of the total energy of the transformed variable trajectory is contained. This is done for each utterance in the training set and the mean  $f_c^j$  (along with standard deviation (SD)) over all utterances is calculated for  $\alpha = 90, 95$ . It is found that the range of  $f_c^j$  for most of the transformed variables  $\tilde{x}_n^j$  is similar to what was observed in the frequency analysis of the articulator position variables in  $x_n^j$  [2]. This is expected since  $\tilde{x}_n^j$  is a linear combination of the articulator position variables  $\mathbf{x}_n$ . This analysis will help us choose the cut-off frequencies while designing filters  $h_n^j$ ,  $\forall j$  in GSC.

### 5.3. Inversion using transformed articulatory features

Since we do not know the true covariance matrix  $\mathbf{K}_{\mathbf{xx}}$  of  $\mathbf{x}_n$ , we estimate  $\mathbf{K}_{\mathbf{xx}}$  using the realizations of  $\mathbf{x}_n$  in the training set as follows:

$$\mathbf{K}_{\mathbf{xx}} = \frac{1}{T-1} \sum_{n=1}^T (\mathbf{x}_n - \bar{\mathbf{x}}) (\mathbf{x}_n - \bar{\mathbf{x}})^T \quad (11)$$

where  $\bar{\mathbf{x}} = \frac{1}{T} \sum_{n=1}^T \mathbf{x}_n$ . Using eigen-decomposition of estimated  $\mathbf{K}_{\mathbf{xx}}$ , we obtain the eigenvector matrix  $\mathbf{V}$ , which is used to transform  $\mathbf{x}_n$  vectors of training, dev and test set to  $\tilde{\mathbf{x}}_n$  (using Eq. (10)), where variables are uncorrelated. Parallel acoustic vectors and transformed articulatory vectors  $\tilde{\mathbf{x}}_n$  of the training set are used to estimate  $\eta_n^{l,j}$  and  $p_n^l$  in a way similar to that described by Ghosh et al. [2]. The GSC (Eq. (1)) is used to estimate the trajectories of  $\tilde{x}_n^{j*}$ ,  $\forall j = 1, \dots, 14$  separately using the acoustic data of the test set. Finally,  $x_n^{j*}$ ,  $\forall j = 1, \dots, 14$  are obtained by transforming  $\tilde{\mathbf{x}}_n^*$  back using  $\mathbf{x}_n^* = [x_n^{1*} \ x_n^{2*} \ \dots \ x_n^{14*}]^T = \mathbf{V} \tilde{\mathbf{x}}_n^* = \mathbf{V} [\tilde{x}_n^{1*} \ \tilde{x}_n^{2*} \ \dots \ \tilde{x}_n^{14*}]^T$ . By this reverse transformation, we correlate different variables so that correlation among them is similar to the inter-articulator correlation observed in the training data. Note that the transformation matrix ( $\mathbf{V}$ ) is learned on the training set. Therefore, it is assumed that correlations between different articulator positions in the test set are similar to those in the training set.

### 5.4. Articulatory inversion results using modified GSC

The proposed approach of utilizing inter-articulator correlation for acoustic-to-articulatory inversion is evaluated separately for

the male and female subjects data in the MOCHA-TIMIT corpus [11]. The accuracy of the inversion is evaluated separately on the test set for both subjects in terms of both root mean squared (RMS) error and Pearson correlation co-efficient [13] between the actual articulatory position in the test set and the position estimated by GSC. The RMS error  $\mathcal{E}$  reflects the average closeness between actual and estimated articulator trajectories. The correlation  $\rho$  indicates how similar the actual and estimated articulator trajectories are.

The dev set is used to tune the cut-off frequency  $f_c^j$  of filter  $h_n^j$  and the trade-off parameter  $C_j$ . For our experiment we considered  $L=200$ . Increasing  $L$  further did not improve the result. For designing articulator specific high-pass filters  $h_n^j$ , we considered an IIR high pass filter with cut-off frequency  $f_c^j$  and stop-band ripple 40 dB down compared to the pass-band ripple similar to that used by Ghosh et al [2]. From Section 5.2, we observed that most of the energy of the transformed variables is below 9-10 Hz and, hence, we choose values of  $f_c^j$  from the following set  $\{f_c^j\} = \left\{1.5 + \frac{(k-1)}{19} (7.5); k = 1, \dots, 20\right\}$ . Similarly, the set of values for  $C_j$  was chosen from the set  $\{.001, .005, .01, 0.05, .1, .5, 1, 5, 10, 50, 100\}$ . The  $f_c^j$  and  $C_j$  combination which yielded the minimum value of the averaged  $\mathcal{E}$  (averaged over all utterances of the dev set) are selected separately for each subject. These best choices of  $f_c^j$  and  $C_j$  are used to perform inversion on the test set using modified GSC. The  $\mathcal{E}$  and  $\rho$  are computed between the actual trajectories and the estimated trajectories for every utterance in the test set. The mean  $\mathcal{E}$  and  $\rho$  (averaged over all utterances in the test set) along with the standard deviation (SD) are shown in Fig. 2 for both the female and male subjects. For comparison, Fig. 2 also shows  $\mathcal{E}$  and  $\rho$  when GSC is directly used in the articulator position variable domain [2].

From Fig. 2, it is clear that for most of the cases, the accuracy of estimates is similar or higher (i.e., lower  $\mathcal{E}$  or higher  $\rho$ ) when inter-articulator correlation is utilized using transformation of variables, but there are cases when GSC on transformed domain either increases the mean  $\mathcal{E}$  or decreases the mean  $\rho$  (e.g. ul<sub>x</sub>, ll<sub>x</sub>, ul<sub>y</sub> for male subject and tt<sub>x</sub> for female subject). However, considering the SD of  $\mathcal{E}$  and  $\rho$ , the changes in the accuracy in terms of  $\mathcal{E}$  and  $\rho$  are insignificant. Thus the gain in performance because of explicitly using inter-articulator correlation by transformation of variable approach is not significant. This supports the conclusions based on the theoretical analysis (Sec. 4) that the correlation among articulators are approximately preserved in GSC and, hence, there is no further benefit by imposing inter-articulator correlation explicitly.

## 6. Conclusions

The analysis of inter-articulator correlation in this paper reveals that acoustic-to-articulatory inversion using GSC approximately preserves inter-articulator correlations although articulatory inversion in GSC is performed for each articulatory trajectory separately. Using both theoretical and experimental analysis, we observe that the smoothness constraints for different articulators are similar and this could be the reason for GSC to approximately preserve correlation in articulatory inversion in practice.

## 7. References

- [1] C. Qin and M. A. Carreira-Perpinan, "An empirical investigation of the nonuniqueness in the acoustic-to-articulatory mapping," *Proc. Interspeech*, pp. 74–77, 2007.
- [2] P. K. Ghosh and S. Narayanan, "A generalized smoothness criterion for acoustic-to-articulatory inversion," *J. Acoust. Soc. Am.*, vol. 128, no. 4, pp. 2162–2172, 2010.
- [3] V. N. Sorokin, A. Leonov, and A. V. Trushkin, "Estimation of sta-

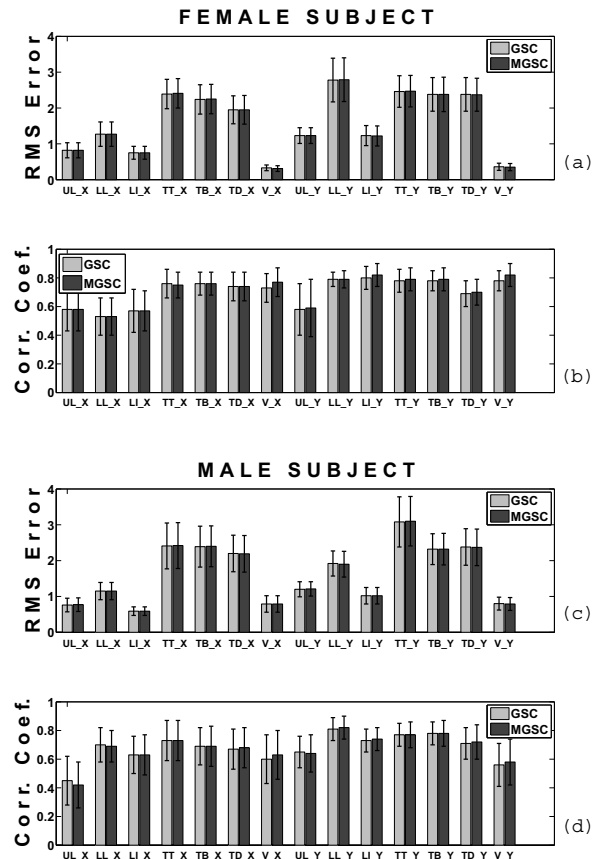


Figure 2: (a)-(b): The average RMS Error ( $\mathcal{E}$ ) and Corr. Coef. ( $\rho$ ) of inversion accuracy using GSC [2] and modified GSC (MGSC) on the test set for the female subject; (c)-(d): (a)-(b) repeated for the male subject. Error bars indicate SD.

bility and accuracy of inverse problem solution for the vocal tract," *Speech Communication*, vol. 30, pp. 55–74, 2000.

- [4] J. Schroeter and M. M. Sondhi, "Dynamic programming search of articulatory codebooks," *Proceedings ICASSP, Glasgow, UK*, pp. 588–591, 1989.
- [5] M. G. Rahim, W. B. Kleijn, J. Schroeter, and C. C. Goodyear, "Acoustic-to-articulatory parameter mapping using an assembly of neural networks," *Proc. ICASSP*, pp. 485–488, 1991.
- [6] J. Hogden, A. Lofqvist, V. Gracco, I. Zlokarnik, P. Rubin, and E. Saltzman, "Accurate recovery of articulator positions from acoustics: New conclusions based on human data," *J. Acoust. Soc. Am.*, vol. 100, no. 3, pp. 1819–1834, 1996.
- [7] T. Toda, A. Black, and K. Tokuda, "Acoustic-to-articulatory inversion mapping with gaussian mixture model," *Proc. ICSLP, Jeju Island, Korea*, pp. 1129–1132, 2004.
- [8] K. Richmond, *Estimating articulatory parameters from the acoustic speech signal*. Ph.D. Thesis, The Centre for Speech Technology Research, Edinburgh University, 2002.
- [9] —, "Mixture density networks, human articulatory data and acoustic-to-articulatory inversion of continuous speech," *Proceedings Workshop on Innovation in Speech Processing WISP*, 2001.
- [10] V. D. Singampalli and P. J. B. Jackson, "Statistical identification of articulation constraints in the production of speech," *Speech Communication*, vol. 51, no. 8, pp. 695–710, 2009.
- [11] A. A. Wrench and K. Richmond, "Continuous speech recognition using articulatory data," *Proc. ICSLP, Beijing, China*, pp. 145–148, 2000.
- [12] H. Stark and J. W. Woods, *Probability and Random Processes with Applications to Signal Processing*, 3rd ed. Prentice Hall, August 3, 2001.
- [13] D. R. Cox and D. V. Hinkley, *Theoretical Statistics*. Chapman & Hall, 1974.