# Voice Conversion Based On Maximum-Likelihood Estimation of Speech Parameter Trajectory

Tomoki Toda, Alan W. Black and Keiichi Tokuda

By:-

Mayank Sirotiya (Y8104036)

Vipul Arora (Y5508)

Under Guidance of:-

Dr. R. Hegde

# Introduction:

This paper describes a method for Spectral Estimation for Voice Conversion application, based on Maximum Likelihood Estimation technique. A Gaussian mixture model (GMM) of the joint probability density of source and target features is employed for performing spectral conversion between speakers. The conventional method converts spectral parameters frame by frame based on the minimum mean square error. However there are some problems associated with this method. In order to address those problems, we propose a conversion method based on the maximum-likelihood estimation of a spectral parameter trajectory.

Voice conversion technology enables a user to transform one person's speech pattern into another pattern with distinct characteristics .A mapping function is used which consists of utterance pairs of source and target voices

# Applications

1. **Speaker conversion** - This technique can modify nonlinguistic information such as voice characteristics while keeping linguistic information unchanged. The resulting mapping function allows the conversion of any sample of the source into that of the target without any linguistic features such as phoneme transcription.

2. **Cross-language speaker conversion** - This conversion framework can straightforwardly be extended to cross-language speaker conversion,  to realize target speakers' voices in various languages by applying the mapping function trained in a certain language into the conversion process in another language.

3. **Modeling of speech production** - Models speech as a combination of a sound source, such as the vocal cords, and a linear acoustic filter, the vocal tract (and radiation characteristic).

4. Narrow-band to wide-band speech for telecommunication

5. Speaking aid, etc.

# Classical approaches and their limitations

A. **Codebook mapping based on hard clustering and discrete mapping** - The converted feature vector at frame 't'is determined by quantizing the source feature vector to the nearest centroid vector of the source codebook and substituting it with a corresponding centroid vector of the mapping codebook as follows:

$$\hat{\boldsymbol{y}}_t = \boldsymbol{c}_m^{(y)}$$

Where, $\hat{\boldsymbol{y}}_t$ is the converted feature vector, and $\boldsymbol{c}_m^{(y)}$ is the centroid vector.

Since this method is based on hard clustering, therefore, large quantization error occurs, which is removed by fuzzy vector quantization.

B. **Fuzzy vector quantization, for soft clustering** - Continuous weights for individual clusters are determined at each frame according to the source feature vector. The converted feature vector is defined as a weighted sum of the centroid vectors of the mapping codebook as follows:

$$\hat{\boldsymbol{y}}_t = \sum_{m=1}^{M} w_{m,t}^{(x)} \boldsymbol{c}_m^{(y)}$$

where M is the number of centroid vectors

C. More variable representations of the converted feature vector are achieved by modeling a difference vector between the source and target feature vectors as follows:

$$\hat{y}_t = x_t + \sum_{m=1}^{M} w_{m,t}^{(x)} \left( c_m^{(y)} - c_m^{(x)} \right)$$

In this method, a very strong correlation between those two vectors is assumed.

D. **Linear multivariate regression (LMR)** - This is based on continuous mapping based on hard clustering

$$\hat{y}_t = A_m x_t + b_m$$

where $A_m$ and $b_m$ are regression parameters.

E. **Gaussian mixture model** - This realizes continuous mapping based on soft clustering as follows:

$$\hat{y}_t = \sum_{m=1}^{M} w_{m,t}^{(x)} \left( A_m x_t + b_m \right)$$

This mapping method is reasonably effective.

# CONVENTIONAL GMM-BASED MAPPING

***Probability Density Function -*** The joint probability density of the source and target feature vectors is

$$P \left( z_t | \lambda^{(z)} \right) = \sum_{m=1}^{M} w_m N \left( z_t; \mu_m^{(Z)}, \Sigma_m^{(z)} \right)$$

where $z_t$ is a joint vector $\left[ x_t^T, y_t^T \right]^T$
and the mean vector and covariance matrix are written as

$$\mu_m^{(Z)} = \begin{bmatrix} \mu_m^{(x)} \\ \mu_m^{(y)} \end{bmatrix}, \Sigma_m^{(z)} = \begin{bmatrix} \Sigma_m^{(xx)} & \Sigma_m^{(xy)} \\ \Sigma_m^{(yx)} & \Sigma_m^{(yy)} \end{bmatrix}$$

***Mapping function -*** Conditional probability density can also be represented as

$$P\left(\mathbf{y}_t|\mathbf{x}_t, \lambda^{(z)}\right) = \sum_{m=1}^{M} P\left(m|\mathbf{x}_t, \lambda^{(z)}\right) P\left(\mathbf{y}_t|\mathbf{x}_t, m, \lambda^{(z)}\right)$$

where

$$P\left(m|\mathbf{x}_t, \lambda^{(z)}\right) = \frac{w_m N\left(\mathbf{x}_t; \mu_m^{(x)}, \Sigma_m^{(xx)}\right)}{\sum_{m=1}^{M} w_n N\left(\mathbf{x}_t; \mu_n^{(x)}, \Sigma_n^{(xx)}\right)}$$

$$P\left(\mathbf{y}_t|\mathbf{x}_t, m, \lambda^{(z)}\right) = N\left(\mathbf{y}_t; \mathrm{E}_{m,t}^{(y)}, \mathrm{D}_m^{(y)}\right)$$

The mean vector and the covariance matrix of m$^{th}$ conditional probability distribution are written as

$$\mathrm{E}_{m,t}^{(y)} = \mu_m^{(y)} + \Sigma_m^{(yx)}\Sigma_m^{(xx)-1}\left(\mathbf{x}_t - \mu_m^{(x)}\right)$$

$$\mathrm{D}_m^{(y)} = \Sigma_m^{(yy)} - \Sigma_m^{(yx)}\Sigma_m^{(xx)-1}\Sigma_m^{(xy)}$$

In conventional method the conversion is based on MMSE as follows

$$\hat{\boldsymbol{y}}_t = E[\boldsymbol{y}_t|\boldsymbol{x}_t]$$
$$= \int P\left(\boldsymbol{y}_t|\boldsymbol{x}_t, \boldsymbol{\lambda}^{(z)}\right)\boldsymbol{y}_t d\boldsymbol{y}_t$$
$$= \int \sum_{m=1}^{M} P\left(m|\boldsymbol{x}_t, \boldsymbol{\lambda}^{(z)}\right) P\left(\boldsymbol{y}_t|\boldsymbol{x}_t, m, \boldsymbol{\lambda}^{(z)}\right)\boldsymbol{y}_t d\boldsymbol{y}_t$$
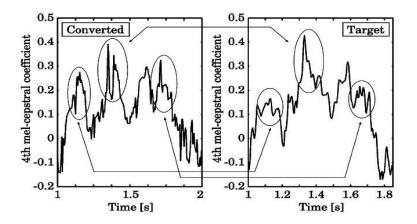$$= \sum_{m=1}^{M} P\left(m|\boldsymbol{x}_t, \boldsymbol{\lambda}^{(z)}\right)\boldsymbol{E}_{m,t}^{(y)}$$

Note that this mapping function has the same form as in our general equation

with $P(m|\boldsymbol{x}_t, \boldsymbol{\lambda}^{(z)}) = w_{m,t}^{(x)}, \Sigma_m^{(yx)}\Sigma_m^{(xx)-1} = \boldsymbol{A}_m,$ and $\boldsymbol{\mu}_m^{(y)} - \Sigma_m^{(yx)}\Sigma_m^{(xx)-1}\boldsymbol{\mu}_m^{(x)} = \boldsymbol{b}_m.$
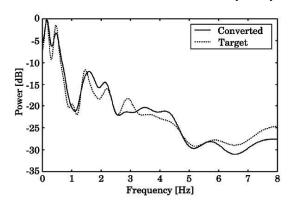
In each mixture component, the conditional target mean vector for the given source feature vector is calculated by a simple linear conversion based on the correlation between the source and target feature vectors problems

# Limitations

1. **Time-independent mapping** - Fig. 2 shows an example of the parameter trajectory converted using the GMM-based mapping function and the natural target trajectory. Although these two trajectories seem similar, they sometimes have different local patterns. Such differences are often observed because the correlation of the target feature vectors between frames is ignored in the conventional mapping.

2. **Oversmoothing** - Fig. 3 shows an example of the converted and natural target spectra. We can see that the converted spectrum is excessively smoothed compared with the natural one. This smoothing undoubtedly causes error reduction of the spectral conversion. However, it also causes quality degradation of the converted speech because the removed structures are still necessary for synthesizing high-quality speech.



# PROPOSED SPECTRAL CONVERSION:

The paper proposes a trajectory based spectral conversion, instead of conventional frame based one. The estimated static feature vector trajectory is given as,

$$\hat{y} = f(x)$$

where,

$$x = \left[ x_1^\top, x_2^\top, \ldots, x_t^\top, \ldots, x_T^\top \right]^\top$$

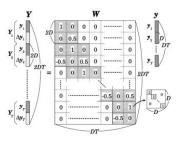$$y = \left[ y_1^\top, y_2^\top, \ldots, y_t^\top, \ldots, y_T^\top \right]^\top$$

## Conversion Considering Dynamic Features:

This method takes into account the feature correlation between frames.

The source and target feature vectors are 2D dimensional

$$X = \left[ X_1^\top, X_2^\top, \ldots, X_t^\top, \ldots, X_T^\top \right]^\top$$

$$Y = \left[ Y_1^\top, Y_2^\top, \ldots, Y_t^\top, \ldots, Y_T^\top \right]^\top$$

where,

$$X_t = \left[ \boldsymbol{x}_t^\top, \Delta\boldsymbol{x}_t^\top \right]^\top \text{ and } Y_t = \left[ \boldsymbol{y}_t^\top, \Delta\boldsymbol{y}_t^\top \right]^\top$$

The static feature vector sequence and static & dynamic feature vector sequence relate as,

$$Y = Wy$$

# Maximum Likelihood Estimation of Parameter Trajectory:

The joint vector is written as,

$$Z_t = \left[ X_t^\top, Y_t^\top \right]^\top$$

The Gaussian Mixture Model (GMM) for joint probability density is trained using conventional training framework.

$$P(Z_t|\boldsymbol{\lambda}^{(Z)})$$

The likelihood function to be maximized is:

$$P\left(Y|X, \boldsymbol{\lambda}^{(Z)}\right) = \prod_{t=1}^{T} \sum_{m=1}^{M} P\left(m|X_t, \boldsymbol{\lambda}^{(Z)}\right) \times P\left(Y_t|X_t, m, \boldsymbol{\lambda}^{(Z)}\right)$$

where,

$$P\left(m|X_t, \boldsymbol{\lambda}^{(Z)}\right) = \frac{w_m \mathcal{N}\left(X_t; \boldsymbol{\mu}_m^{(X)}, \boldsymbol{\Sigma}_m^{(XX)}\right)}{\sum\limits_{n=1}^{M} w_n \mathcal{N}\left(X_t; \boldsymbol{\mu}_n^{(X)}, \boldsymbol{\Sigma}_n^{(XX)}\right)}$$

$$P\left(Y_t|X_t, m, \boldsymbol{\lambda}^{(Z)}\right) = \mathcal{N}\left(Y_t; E_{m,t}^{(Y)}, D_m^{(Y)}\right)$$

here,

$$E_{m,t}^{(Y)} = \boldsymbol{\mu}_m^{(Y)} + \boldsymbol{\Sigma}_m^{(YX)} \boldsymbol{\Sigma}_m^{(XX)-1} \left(X_t - \boldsymbol{\mu}_m^{(X)}\right)$$

$$D_m^{(Y)} = \boldsymbol{\Sigma}_m^{(YY)} - \boldsymbol{\Sigma}_m^{(YX)} \boldsymbol{\Sigma}_m^{(XX)-1} \boldsymbol{\Sigma}_m^{(XY)}$$

This relation is derived as follows [2].

## Derivation of Conditional Probability:

The joint vector with Gaussian distribution $\mathcal{N}(\mathbf{x}|\mu, \Sigma)$ is written as,

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$$

with mean and covariance matrices as,

$$\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \qquad \Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

For convenience, inverse of covariance matrix, known as Precision matrix, is used,

$$\Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

As inverse of a symmetric matrix is also symmetric, hence, $\Lambda_{aa}$ and $\Lambda_{bb}$ are symmetric, and $\Lambda_{aa} = \Lambda_{bb}^T$.

From the product rule of probability, we see that the conditional distribution can be evaluated from the joint distribution $\mathbf{p(x)} = \mathbf{p(x_a, x_b)}$ simply by fixing $\mathbf{x_b}$ to the observed value and normalizing the resulting expression to obtain a valid probability distribution over $\mathbf{x_a}$. Instead of performing this normalization explicitly, we can obtain the solution more efficiently by considering the quadratic form,

$$-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})=$$

$$-\frac{1}{2}(\mathbf{x}_a-\boldsymbol{\mu}_a)^{\mathrm{T}}\boldsymbol{\Lambda}_{aa}(\mathbf{x}_a-\boldsymbol{\mu}_a)-\frac{1}{2}(\mathbf{x}_a-\boldsymbol{\mu}_a)^{\mathrm{T}}\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b-\boldsymbol{\mu}_b)$$

$$-\frac{1}{2}(\mathbf{x}_b-\boldsymbol{\mu}_b)^{\mathrm{T}}\boldsymbol{\Lambda}_{ba}(\mathbf{x}_a-\boldsymbol{\mu}_a)-\frac{1}{2}(\mathbf{x}_b-\boldsymbol{\mu}_b)^{\mathrm{T}}\boldsymbol{\Lambda}_{bb}(\mathbf{x}_b-\boldsymbol{\mu}_b)$$

Comparing this with the general Gaussian exponent form,

$$-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})=-\frac{1}{2}\mathbf{x}^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\mathbf{x}+\mathbf{x}^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu}+\text{const}$$

we obtain, that the second order term is,

$$-\frac{1}{2}\mathbf{x}_a^{\mathrm{T}}\boldsymbol{\Lambda}_{aa}\mathbf{x}_a$$

and the first order term is,

$$\mathbf{x}_a^{\mathrm{T}}\left\{\boldsymbol{\Lambda}_{aa}\boldsymbol{\mu}_a-\boldsymbol{\Lambda}_{ab}(\mathbf{x}_b-\boldsymbol{\mu}_b)\right\}$$

Hence, we obtain the conditional mean and variance as,

$$\boldsymbol{\mu}_{a|b}=\boldsymbol{\mu}_a+\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}(\mathbf{x}_b-\boldsymbol{\mu}_b)$$
$$\boldsymbol{\Sigma}_{a|b}=\boldsymbol{\Sigma}_{aa}-\boldsymbol{\Sigma}_{ab}\boldsymbol{\Sigma}_{bb}^{-1}\boldsymbol{\Sigma}_{ba}.$$

## EM Algorithm:

$$\hat{\boldsymbol{y}}=\arg\max P\left(\boldsymbol{Y}|\boldsymbol{X},\boldsymbol{\lambda}^{(Z)}\right)$$

This direct maximization is very complex, hence EM algorithm is used, in which the following auxiliary function is iteratively maximized with respect to $\hat{\boldsymbol{y}}$,

$$Q(\boldsymbol{Y},\hat{\boldsymbol{Y}})=\sum_{\text{all }\boldsymbol{m}}P\left(\boldsymbol{m}|\boldsymbol{X},\boldsymbol{Y},\boldsymbol{\lambda}^{(Z)}\right)\log P\left(\hat{\boldsymbol{Y}},\boldsymbol{m}|\boldsymbol{X},\boldsymbol{\lambda}^{(Z)}\right)$$

$$=\sum_{t=1}^{T}\sum_{m=1}^{M}P\left(m|\boldsymbol{X}_t,\boldsymbol{Y}_t,\boldsymbol{\lambda}^{(Z)}\right)\log P\left(\hat{\boldsymbol{Y}}_t,m|\boldsymbol{X}_t,\boldsymbol{\lambda}^{(Z)}\right)$$

$$=\sum_{t=1}^{T}\sum_{m=1}^{M}\gamma_{m,t}\left(-\frac{1}{2}\hat{\boldsymbol{Y}}_t^{\top}\boldsymbol{D}_m^{(Y)-1}\hat{\boldsymbol{Y}}_t+\hat{\boldsymbol{Y}}_t^{\top}\boldsymbol{D}_m^{(Y)-1}\boldsymbol{E}_{m,t}^{(Y)}\right)+\overline{K}$$

$$=\sum_{t=1}^{T}-\frac{1}{2}\hat{\boldsymbol{Y}}_t^{\top}\overline{\boldsymbol{D}_t^{(Y)-1}}\hat{\boldsymbol{Y}}_t+\hat{\boldsymbol{Y}}_t^{\top}\overline{\boldsymbol{D}_t^{(Y)-1}\boldsymbol{E}_t^{(Y)}}+\overline{K}$$

$$=-\frac{1}{2}\hat{\boldsymbol{Y}}^{\top}\overline{\boldsymbol{D}^{(Y)-1}}\hat{\boldsymbol{Y}}+\hat{\boldsymbol{Y}}^{\top}\overline{\boldsymbol{D}^{(Y)-1}\boldsymbol{E}^{(Y)}}+\overline{K}$$

$$=-\frac{1}{2}\hat{\boldsymbol{y}}^{\top}\boldsymbol{W}^{\top}\overline{\boldsymbol{D}^{(Y)-1}}\boldsymbol{W}\hat{\boldsymbol{y}}+\hat{\boldsymbol{y}}^{\top}\boldsymbol{W}^{\top}\overline{\boldsymbol{D}^{(Y)-1}\boldsymbol{E}^{(Y)}}+\overline{K}$$

where,

$$\overline{\boldsymbol{D}^{(Y)-1}}=\text{diag}\left[\overline{\boldsymbol{D}_1^{(Y)-1}},\overline{\boldsymbol{D}_2^{(Y)-1}},\dots,\right.$$
$$\left.\overline{\boldsymbol{D}_t^{(Y)-1}},\dots\overline{\boldsymbol{D}_T^{(Y)-1}}\right]$$

$$\overline{\boldsymbol{D}^{(Y)-1}\boldsymbol{E}^{(Y)}}=\left[\overline{\boldsymbol{D}_1^{(Y)-1}\boldsymbol{E}_1^{(Y)}}^{\top},\overline{\boldsymbol{D}_2^{(Y)-1}\boldsymbol{E}_2^{(Y)}}^{\top},\dots,\right.$$
$$\left.\overline{\boldsymbol{D}_t^{(Y)-1}\boldsymbol{E}_t^{(Y)}}^{\top},\dots,\overline{\boldsymbol{D}_T^{(Y)-1}\boldsymbol{E}_T^{(Y)}}^{\top}\right]^{\top}$$

$$\overline{\boldsymbol{D}_t^{(Y)-1}}=\sum_{m=1}^{M}\gamma_{m,t}\boldsymbol{D}_m^{(Y)-1}$$

$$\overline{\boldsymbol{D}_t^{(Y)-1}\boldsymbol{E}_t^{(Y)}}=\sum_{m=1}^{M}\gamma_{m,t}\boldsymbol{D}_m^{(Y)-1}\boldsymbol{E}_{m,t}^{(Y)}$$

$$\gamma_{m,t}=P(m|\boldsymbol{X}_t,\boldsymbol{Y}_t,\boldsymbol{\lambda}^{(Z)}).$$

The derivative of auxiliary function with respect to $\hat{\boldsymbol{y}}$ is,

$$\frac{\partial Q(\boldsymbol{Y},\hat{\boldsymbol{Y}})}{\partial y}=-\boldsymbol{W}^{\top}\overline{\boldsymbol{D}_{\hat{m}}^{(Y)-1}}\boldsymbol{W}\boldsymbol{y}+\boldsymbol{W}^{\top}\overline{\boldsymbol{D}_{\hat{m}}^{(Y)-1}}\boldsymbol{E}_{\hat{m}}^{(Y)}$$

Equating it to zero, we get,

$$\hat{y} = \left(W^\top \overline{D^{(Y)^{-1}}} W\right)^{-1} W^\top \overline{D^{(Y)^{-1}} E^{(Y)}}$$

In the expectation step, we calculate $\gamma_{m,t}$ and in the maximization step, we estimate $\hat{y}$.

## Approximation with suboptimum mixture sequence:

The log likelihood function is approximated with a single mixture component sequence as

$$\mathcal{L} = \log P\left(\hat{m}|X, \lambda^{(Z)}\right) P\left(Y|X, \hat{m}, \lambda^{(Z)}\right)$$

where,

$$\hat{m} = \operatorname{argmax} P\left(m|X, \lambda^{(Z)}\right)$$

This log likelihood function is easy to maximize by direct differentiation. Hence, we obtain,

$$\hat{y} = \left(W^\top D_{\hat{m}}^{(Y)^{-1}} W\right)^{-1} W^\top D_{\hat{m}}^{(Y)^{-1}} E_{\hat{m}}^{(Y)}$$
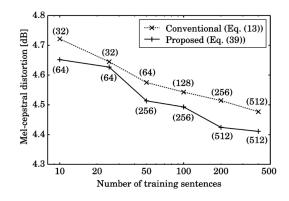
where,

$$E_{\hat{m}}^{(Y)} = \left[E_{\hat{m}_1,1}^{(Y)}, E_{\hat{m}_2,2}^{(Y)}, \ldots, \right.$$
$$\left. E_{\hat{m}_t,t}^{(Y)}, \ldots, E_{\hat{m}_T,T}^{(Y)}\right]$$
$$D_{\hat{m}}^{(Y)^{-1}} = \operatorname{diag}\left[D_{\hat{m}_1}^{(Y)^{-1}}, D_{\hat{m}_2}^{(Y)^{-1}}, \ldots, \right.$$
$$\left. D_{\hat{m}_t}^{(Y)^{-1}}, \ldots, D_{\hat{m}_T}^{(Y)^{-1}}\right]$$

The approximated solution reduces the computational cost. The preliminary experiments demonstrated that there were no large differences in the results obtained by EM algorithm and the approximated method.

# Results:

The evaluation using mel-cepstral distortions show that the proposed method is quite better than the conventional method. This is because the proposed method takes into account the interframe correlation for parameter trajectory estimation, whereas conventional method ignores it.



# Summary:

In this paper, we mainly studied following concepts:
- GMM based Feature mapping
- Conditional Gaussian Distributions
- Maximum Likelihood technique for GMMs
- Expectation Maximization Algorithm