ADAPTATION OF A PREDICTIVE MODEL OF TONGUE SHAPES Chao Qin and Miguel Á. Carreira-Perpiñán EECS, School of Engineering, University of California, Merced, USA

Abstract

It is possible to recover the full midsagittal contour of the tongue with submillimetric accuracy from the location of just 3–4 landmarks on it. This involves fitting a predictive mapping from the landmarks to the contour using a training set consisting of contours extracted from ultrasound recordings. However, extracting sufficient contours is a slow and costly process. Here, we consider adapting a predictive mapping obtained for one condition (such as a given recording session, recording modality, speaker or speaking style) to a new condition, given only a few new contours and no correspondences. We propose an extremely fast method based on estimating a 2D-wise linear alignment mapping, and show it recovers very accurate predictive models from about 10 new contours.

Motivation and idea

- Extracting sufficient good contours to estimate a predictive model is hard
- Need to obtain models for new conditions (e.g. a new recording session, recording modality, speaker or speaking style) \Rightarrow adaptation
- The adaptation problem: Given the predictive mapping f (that maps landmarks to contours) and a small adaptation set of N contours $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$, want to adapt f by estimating an invertible mapping g (with few parameters) that maps new data (from the new speaker) to old data (from the reference speaker)
- Advantages of our adaptation method:
- No need for any correspondence
- Little data needed for adaptation
- Few parameters to estimate \Rightarrow extremely fast
- Applicable to both linear and nonlinear predictive mappings
- Extendable to 3D shapes

Idea of the method

- Learn a predictive mapping f from a dataset $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$ containing N contours
- **2** Estimate 2D-wise transformation $\mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$
- ${f eta}$ Obtain the adapted mapping ${f f}' = {f g}_{f v}^{-1} \circ {f f} \circ {f g}_{f x}$



Predictive model of tongue shapes

The prediction problem: given the 2D locations of K landmarks located on the tongue midsagittal contour (x), reconstruct the entire contour (y), represented by P 2D points

- Linear mapping: f(x) = Wx + w
- Radial basis function (RBF) network: $f(x) = W\Phi(x) + w$ with M Gaussian basis functions $\phi_m(\mathbf{x}) = \exp\left(-\frac{1}{2} \|(\mathbf{x} - \boldsymbol{\mu}_m) / \sigma\|^2\right)$ and width σ

Adaptation based on feature normalization

• Key aspect: apply the same transformation g to each 2D point of an x- or y-contour.

$$ilde{\mathbf{x}} = \mathbf{g}_{\mathbf{x}}(\mathbf{x}) = egin{pmatrix} \mathbf{A} oldsymbol{x}_1 + \mathbf{b} \\ \dots \\ \mathbf{A} oldsymbol{x}_K + \mathbf{b} \end{pmatrix} \qquad ilde{\mathbf{y}} = \mathbf{g}_{\mathbf{y}}(\mathbf{y}) = egin{pmatrix} \mathbf{A} oldsymbol{y}_1 + \mathbf{b} \\ \dots \\ \mathbf{A} oldsymbol{y}_P + \mathbf{b} \end{pmatrix}$$

- The adapted predictive mapping f' is given by $\mathbf{g}_{\mathbf{v}}^{-1} \circ \mathbf{f} \circ \mathbf{g}_{\mathbf{x}}$
- Advantage of 2D-wise alignment mapping $\mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$: -Linear \Rightarrow invertible and 6 parameters to estimate, $A_{2\times 2}$ and $b_{2\times 1}$.
- -Requires very little adaptation data (one contour is enough if P > 2 points)
- To estimate g, we minimize the error function

$$F(\mathbf{A}, \mathbf{b}) = \sum_{n=1}^{N} ||\mathbf{g}_{\mathbf{y}}(\mathbf{y}_n) - \mathbf{f}(\mathbf{g}_{\mathbf{x}}(\mathbf{x}_n))|$$

which is easier to minimize than the square error error

$$E(\mathbf{A}, \mathbf{b}) = \sum_{n=1}^{N} \|\mathbf{y}_n - \mathbf{g}_{\mathbf{y}}^{-1}(\mathbf{f}(\mathbf{g}_{\mathbf{x}}(\mathbf{x}_n)))$$

• Solution:

-Linear mapping: unique solution from a positive definite 6×6 linear system -RBF: we apply the BFGS algorithm and find no local optima. In practice, superlinear convergence in ~ 10 iterations

Computational complexity

Theory	
N = 10, M = 500, P = 24, K =	= 3

Linear mapping	
$\mathcal{O}(32NP)$	$\mathcal{O}(14NM(P +$
4 ms	

Conclusion

We have proposed a fast (< 1 second) adaptation method for RBF predictive mappings of tongue contours and shown that a few contours suffice to achieve near-optimal accuracy. In synthetic transformations and in the problem of correcting for misalignments between different ultrasound recording sessions, just one contour reduces the error per contour point below 1 mm, and 10-20 contours bring it within 5–10% of the one obtained by training from scratch on a large training set. Future work will involve using more flexible adaptation models (e.g. using a different A and b per 2D point) and testing the model with data from different speakers, when the latter becomes available.

 $f(\mathbf{g}_{\mathbf{x}}(\mathbf{x}))$

 $\mathbf{g}_{\mathbf{x}}(\mathbf{x})$

Work funded by NSF award IIS-0711186

RBF K)) per BFGS iteration 0.1 s

Experimental results

Datasets

- 8 671 ultrasound tongue contours (P = 24) from a set of 22 British TIMIT sentences for one Scottish speaker, maaw0 • Datasets **S1** (3727 contours from 10 utterances) from session 1 and **S2** (4944 contours from 12 utterances) from session 2

Predictive models

- K landmarks were chosen optimally from the P contour points • Linear mapping: as a baseline and provide initial $\{A, b\}$ for the RBF • RBF: M = 500, $\sigma = 55$, and $\lambda = 10^{-4}$, trained by cross-validation on
- the $2\,236$ contours of **S1**

Comparison methods

- Retraining: trains the predictive mapping from scratch on the adaptation data
- PCA alignment: finds $\{A, b\}$ by matching the mean and covariance (principal axes' angle and variance) of the original and the adaptation datasets

• Ground truth: retrains the predictive model with abundant data Adaptation tasks

- Task 1 recover a known transformation $\{A_0, b_0\}$: $A_0 = \begin{pmatrix} 0.26 & 1.82 \\ -0.36 & 0.26 \end{pmatrix}$, $b_0 = \begin{pmatrix} 52.8 \\ 67.4 \end{pmatrix}$
- Task 2 alignment between recording sessions: to adapt the predictive model f of **S1** to data from **S2**





- **S1** (recording session 1) **S2** (recording session 2)

0 - 0 show pointwise RMSE E after adaptation/retraining w.r.t N contours (using K = 3 landmarks) and K landmarks (using N = 10). Errorbars are over 10 random choices of N. Θ shows a joint, corrected dataset (S1,S2), obtained by using all the S2 contours to adapt f from S1 and aligning S2 to S1 with the resulting g. • Adaptation is much better than retraining and PCA alignment especially when the adaptation data is limited • Adaptation is effective to minimize the misalignment among different recording sessions