Trajectory Analysis of Speech using Continuous-State Hidden Markov Models

P. Weber, S. M. Houghton, C. J. Champion, M. J. Russell and P. Jančovič

University of Birmingham, UK

ICASSP, 7 May 2014

Continuously Changing Smooth Trajectories

Speech is not a set of discrete states

- rather a series of smooth, continuous transitions.



"he will allow a rare lie"

HMS Linear Dwell-Transition Model



Holmes, Mattingley, Shearme, 'Speech Synthesis by Rule', (1964) (HMS).

- A sequence of stationary periods linked by smooth transitions.
- Piece-wise linear approximation.
- Dwell (articulator) target frequencies, transitions.

Given outputs generated according to the HMS Model, fit a continuous sequence of trajectories to them ...



- ... and thus recover the sequence of phonemes.
- Continuous-State HMM algorithm.

Speech is modelled as a sequence of constant dwells and linear transitions. Dwell represents target sound (e.g. phoneme):

- Persists for $t_{\phi} \sim D_{\phi}$ time steps.
- Canonical phoneme frequency targets, \mathbf{f}_{ϕ} .
- Actual frequency targets: noisy realisations, $\mathbf{f}_t \sim \mathcal{N}_d(\mathbf{f}_{\phi}, \mathbf{A})$.
- Realisation variance allows for systematic departures from the mean (e.g. speaker dependence if that is not modelled separately).
- Noisy observations, $\mathbf{y}_t \sim \mathcal{N}_d(\mathbf{f}_t, \mathbf{E})$.

Transition represents smooth movement of articulators between dwells:

- Persists for $t_T \sim D_T$ time steps.
- Frequencies transition linearly between realised targets.
- Noisy observations about the linear transition.

Image: A image: A

CS-HMM: Assumptions



Red: canonical targets, Green: realised targets, Blue: observations.

Weber, Houghton *et al.* (Birmingham)

Trajectory Analysis (CS-HMM)

ICASSP, 7 May 2014 6 / 18

Continuous component

- x: realised target frequencies at time t (dwells and transitions),
- s: slopes at time t (transitions),

Discrete component

- identify whether in dwell or transition,
- *h*: time steps in current dwell/transition,
- identity of current phoneme,
- Phonetic history.

Probability information about an infinite set of states,

• Alpha value $\alpha_{t-1}(\mathbf{x})$:

- sum of probabilities over all paths leading to state at $t-1\ldots$
- ... consistent with discrete history...
- ... given path transitions and observations.
- Parametric form: scaled Gaussian,

$$\begin{split} \alpha_{t-1}(\mathbf{x}) &= \mathcal{K}_{t-1} \mathcal{N}_d(\mathbf{x} - \mu, \mathbf{P}), \\ \text{where} \quad \mathcal{N}_d(\mathbf{x}, \mathbf{P}) &= (2\pi)^{-d/2} |\mathbf{P}|^{1/2} \exp\{-\frac{1}{2} \mathbf{x}^T \mathbf{P} \mathbf{x}\}. \end{split}$$

 μ and P are the mean and precision (1/variance) of the distribution over the current continuous state.

 K_t is the sum of probabilities of paths consistent with the hypothesis. d is the dimension of the space.

Learn the parameters in the system, e.g.

- **③** Per-phoneme canonical target frequencies $\mathbf{f}_{\phi} \ (\approx 40 \times 3)$.
- Target frequency co-variance matrix A (6).
- **Observation co-variance matrix** E (6).
- Timing model (can be anything).

CS-HMM: Recovery

Assume dwell start: Initialise one hypothesis per phoneme

$$\alpha_0(\mathbf{x}) = \mathcal{N}_d(\mathbf{x} - \mathbf{f}_\phi, \mathbf{A}).$$

Step through dwell. Observe y_t, assumed drawn from N_d(x, E).
Update hypothesis to take account of observation

$$\alpha_t(\mathbf{x}) = K_{t-1} \mathcal{N}_d(\mathbf{x} - \mu_{t-1}, \mathbf{P}_{t-1}) \mathcal{N}_d(\mathbf{y}_t - \mathbf{x}, \mathbf{E})$$
$$= K_t \mathcal{N}_d(\mathbf{x} - \mu_t, \mathbf{P}_t).$$

where

$$\begin{split} \mathbf{P}_{t} &= \mathbf{P}_{t-1} + \mathbf{E}, \\ \mu_{t} &= \mathbf{P}_{t}^{-1} (\mathbf{P}_{t-1} \mu_{t-1} + \mathbf{E} \mathbf{y}_{t}), \\ \mathbf{K}_{t} &= \mathbf{K}_{t-1} \, \mathcal{N}_{d} \Big(\mathbf{y}_{t} - \mu_{t-1}, \big(\mathbf{P}_{t-1}^{-1} + \mathbf{E}^{-1} \big)^{-1} \Big). \end{split}$$

10 / 18

Similar (but more complicated) formulae allow us to step through a transition and to move between dwells and transitions.

- At each tick during a dwell we have a choice between continuing the dwell and entering a transition.
- At each tick during a transition we have a choice between continuing the transition and entering a dwell for each phoneme in the inventory.
- Hence we branch on hypotheses at each step.
- And to keep their number within bounds, we threshold on K_t .
- For a full development of the mathematics, see the papers on our website (www.birmingham.ac.uk/srbs).

CS-HMM Process Illustration



- o Showing just one trajectory.
- o Only two phonemes in inventory.



- o Concentrating on one hypothesis.
- o Realised target adjusts.





o New hypotheses for each phoneme target. Limited experimentation to illustrate the technique.

- Recover TIMIT phoneme sequences.
- Vocal Tract Resonances [Deng et al., 2006].
- Learn and recognise single utterances.

Recovery of Phoneme Sequences from TIMIT



Recovery of Phoneme Sequences from TIMIT



Utterance test/dr2/mwew0/sx11: 'he will allow a rare lie'.

Transcription:/hh iy w l ah l aw er r eh r l ay/Recovery:/hh iy w l ah l aw ah er r eh r l ay/

Weber, Houghton et al. (Birmingham)

Trajectory Analysis (CS-HMM)

ICASSP, 7 May 2014

15 / 18

Recovery of Phoneme Sequences from TIMIT



Utterance test/dr2/mwew0/sx11: 'he will allow a rare lie'.

Transcription:/hh iy w l ah l aw er r eh r l ay/Recovery:/hh iy w l ah l aw ah er r eh r l ay/

Weber, Houghton et al. (Birmingham)

Trajectory Analysis (CS-HMM)

ICASSP, 7 May 2014

16 / 18

The CS-HMM algorithm is able to successfully recover phoneme sequences in a controlled environment

- using a more faithful model of speech,
- and a very limited parameter set.

Some developments are clear:

- accurate training of the phoneme inventory,
- automated parameter learning,
- handle systematic variation in speech,
- how does the model extend to non-sonorant speech?

Thank you.

Questions?

http://www.birmingham.ac.uk/SRbS/