

Statistical Parametric Speech Synthesis Based on Product of Experts

Heiga Zen & Mark Gales

TOSHIBA

Leading Innovation >>>

Yoshihiko Nankaku & Keiichi Tokuda



Background

Statistical parametric synthesis based on HMMs

- Advantage

- * Flexibility to change voice characteristics
- * Small footprint
- * Robustness

- Limitation

- * Quality

- Major factors of quality degradation

- * Vocoding
- * Accuracy of acoustic models (AMs)
- * Over-smoothing

This work addresses accuracy of AMs & over-smoothing

Combination of Multiple-Level Acoustic Models

Combine multiple AMs to reduce over-smoothing

- * Training; estimate multiple-level AMs *individually*

$$\hat{\lambda}_i = \arg \max_{\lambda_i} p(f_i(\mathbf{c}) \mid \lambda_i) \quad i = 1, \dots, M$$

- * Synthesis; generate \mathbf{c} that *jointly* maximize output probs from AMs

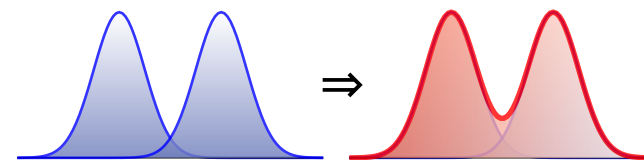
$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} \sum_{i=1}^M \alpha_i \log p(f_i(\mathbf{c}) \mid \hat{\lambda}_i)$$

- * Feature function, $f_i(\mathbf{c})$, extracts acoustic feats for i -th AM from \mathbf{c}
 - e.g., dynamic feats, DCT, average, summation, global variance
- * Parameters of AMs, λ_i , are trained *independently*
 - Use weights to control balance among AMs
- * Weights, α_i , are determined by *held-out data* (or tuned manually)

Mixture Model vs Product Model

Mixture of experts

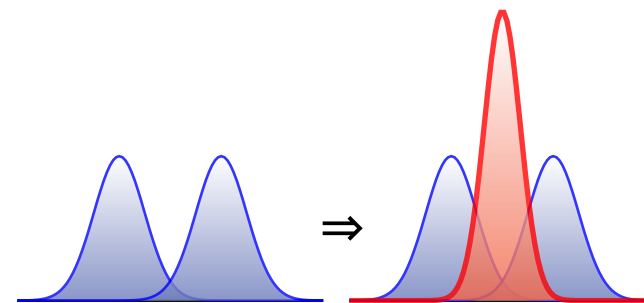
$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \sum_{i=1}^M \alpha_i p(f_i(\mathbf{c}) \mid \lambda_i)$$



- * Data is generated from **union** of experts
- * Robust for modeling data with many variations
- * GMM \rightarrow Mixture of Gaussians

Product of experts [Hinton;'02]

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$



- * Data is generated from **intersection** of experts
- * Efficient for modeling data with many constraints
- * PoG \rightarrow Product of Gaussians

Combination of Multiple-Level AMs as PoE

Combination of multiple AMs can be viewed as PoE

$$\begin{aligned}\hat{\mathbf{c}} &= \arg \max_{\mathbf{c}} p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \arg \max_{\mathbf{c}} \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} \\ &= \arg \max_{\mathbf{c}} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} = \arg \max_{\mathbf{c}} \sum_{i=1}^M \alpha_i \log p(f_i(\mathbf{c}) \mid \lambda_i)\end{aligned}$$

* Generating \mathbf{c} from combination of multiple AMs

→ Equivalent to generating \mathbf{c} from PoE consisting of AMs

* Regarding combination of multiple AMs as PoE

→ **Jointly** estimate multiple AMs

$$\{\hat{\lambda}_1, \dots, \hat{\lambda}_M\} = \arg \max_{\lambda_1, \dots, \lambda_M} \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$

Combination of Multiple-Level AMs as PoE

Combination of multiple AMs can be viewed as PoE

$$\begin{aligned}\hat{\mathbf{c}} &= \arg \max_{\mathbf{c}} p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \arg \max_{\mathbf{c}} \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} \\ &= \arg \max_{\mathbf{c}} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i} = \arg \max_{\mathbf{c}} \sum_{i=1}^M \alpha_i \log p(f_i(\mathbf{c}) \mid \lambda_i)\end{aligned}$$

* Generating \mathbf{c} from combination of multiple AMs

→ Equivalent to generating \mathbf{c} from PoE consisting of AMs

* Regarding combination of multiple AMs as PoE

→ **Jointly** estimate multiple AMs

- AMs become **complementary**

- **held-out data** to estimate weights is no longer required

Product of Gaussians

Product of Gaussians (PoG)

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \mathcal{N}(f_i(\mathbf{c}); \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

- * Special case of PoE; All experts are Gaussian
- * If all feature functions are linear
 - PoG also becomes Gaussian
 - Normalization constant

$$Z = \int \prod_{i=1}^M \mathcal{N}(f_i(\mathbf{c}); \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{c}$$

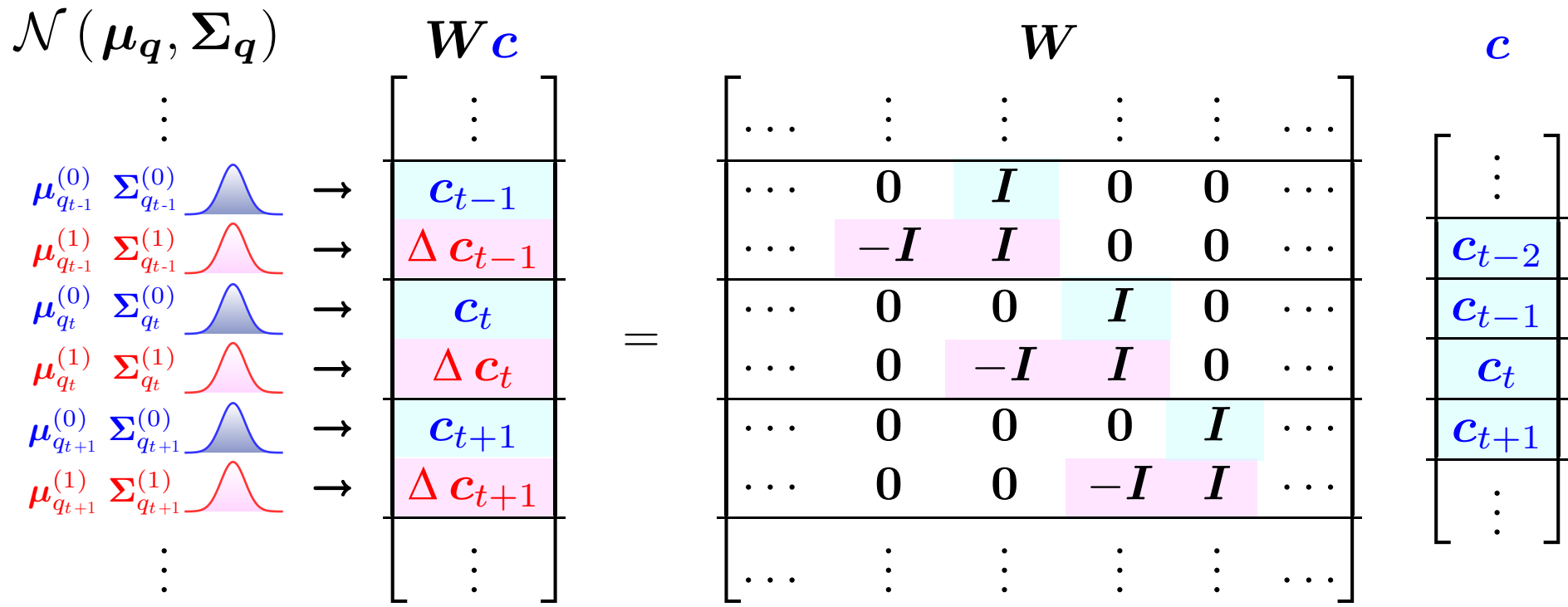
can be computed in **closed form**

Trajectory HMM as Product of Gaussians

Trajectory HMM can be viewed as PoG [Williams;'05, Zen;'07]

$$p(\mathbf{c} | \lambda) = \sum_{\forall \mathbf{q}} p(\mathbf{c} | \mathbf{q}, \lambda) p(\mathbf{q} | \lambda)$$

$$p(\mathbf{c} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$



Trajectory HMM as Product of Gaussians

Trajectory HMM can be viewed as PoG [Williams;'05, Zen;'07]

$$p(\mathbf{c} | \lambda) = \sum_{\forall \mathbf{q}} p(\mathbf{c} | \mathbf{q}, \lambda) p(\mathbf{q} | \lambda)$$

$$p(\mathbf{c} | \mathbf{q}, \lambda) = \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}}) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{W}\mathbf{c}; \boldsymbol{\mu}_{\mathbf{q}}, \boldsymbol{\Sigma}_{\mathbf{q}})$$

$$= \frac{1}{Z_{\mathbf{q}}} \prod_{t=1}^T \prod_{d=0}^2 \mathcal{N}\left(f_t^{(d)}(\mathbf{c}); \boldsymbol{\mu}_{q_t}^{(d)}, \boldsymbol{\Sigma}_{q_t}^{(d)}\right) \quad f_t^{(d)}(\mathbf{c}) : d\text{-th dyn feat at frame } t$$

$$Z_{\mathbf{q}} = \int \prod_{t=1}^T \prod_{d=0}^2 \mathcal{N}\left(f_t^{(d)}(\mathbf{c}); \boldsymbol{\mu}_{q_t}^{(d)}, \boldsymbol{\Sigma}_{q_t}^{(d)}\right) d\mathbf{c}$$

- * Experts are Gaussians, feature functions are dynamic features
- * Gaussian experts are multiplied over time
- * Training/adaptation/search algorithms have been derived

Linear Feature Function with Gaussian Experts

Combining multiple-level AMs as PoE

- * Multiple-level AMs often use linear feature functions w/ Gaussians
 - DCT [Latorre;'08, Qian;'09], average [Wang;'08], sum [Ling;'06, Gao;'08]
- * PoEs become the same form as trajectory HMM
 - Training algorithm for trajectory HMM are applicable

Example: state & phoneme duration models [Ling;'06]

$$\begin{aligned}
 p(\mathbf{d} | \lambda) &= \frac{1}{Z} \mathcal{N}(\mathbf{W}\mathbf{d}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\
 &= \frac{1}{Z} \prod_{i=1}^P \prod_{j=1}^{N_i} \mathcal{N}(d_{ij}; \xi_{ij}, \sigma_{ij}) \\
 &\quad \cdot \prod_{k=1}^P \mathcal{N}(p_k; \nu_k, \omega_k)
 \end{aligned}$$

d_{ij} : duration of state j in phoneme i

	\mathbf{W}						\mathbf{d}
d_{11}	1	0	0				...
d_{12}	0	1	0				...
d_{13}	0	0	1				...
p_1	1	1	1		$\mathbf{0}$...
d_{21}							d_{11}
d_{22}							d_{12}
d_{23}							d_{13}
p_2					$\mathbf{0}$...
							d_{21}
							d_{22}
							d_{23}
\vdots	\vdots	\vdots	\vdots				\vdots
							...

General PoE (Non-Linear/Non-Gaussian)

General form of PoE

$$p(\mathbf{c} \mid \lambda_1, \dots, \lambda_M) = \frac{1}{Z} \prod_{i=1}^M \{p(f_i(\mathbf{c}) \mid \lambda_i)\}^{\alpha_i}$$

- * Feature functions can be non-linear, experts can be non-Gaussian
- * Normalization term has no closed form
- * Training is complex, usually normalization term is approximated

Example: global variance (GV) [Toda;'07]

$$p(\mathbf{c} \mid \mathbf{q}, \lambda, \lambda_{GV}) = \frac{1}{Z_{\mathbf{q}}} \mathcal{N}(\mathbf{c}; \bar{\mathbf{c}}_{\mathbf{q}}, \mathbf{P}_{\mathbf{q}})^{\alpha} \mathcal{N}(f_v(\mathbf{c}); \boldsymbol{\mu}_v, \boldsymbol{\Sigma}_v)$$

$$f_v(\mathbf{c}) = \frac{1}{T} \sum_{t=1}^T \text{diag} \left[(\mathbf{c}_t - \bar{\mathbf{c}}) (\mathbf{c}_t - \bar{\mathbf{c}})^{\top} \right] : \text{intra-utt variance, **quadratic**}$$

Contrastive Divergence Learning

Contrastive divergence learning [Hinton;'02]

- * Training algorithm for general PoE
- * Combination of sampling & gradient methods

1. Draw J samples from PoE

$$\mathbf{c}^{(j)} \sim p(\mathbf{c} | \boldsymbol{\lambda}) \quad j = 1, \dots, J \quad \boldsymbol{\lambda} = \{\lambda_1, \dots, \lambda_M\}: \text{PoE model params}$$

2. Compute approximated derivative of log likelihood w.r.t. $\boldsymbol{\lambda}$

$$\frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \approx \underbrace{\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_0}}_{\text{expectation over data}} - \underbrace{\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_J}}_{\text{expectation over samples}}$$

3. Update model params using gradient method

$$\boldsymbol{\lambda}' = \boldsymbol{\lambda} - \eta \cdot \left(\left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_0} - \left\langle \frac{\partial \log p(\mathbf{c} | \boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} \right\rangle_{p_j} \right) \quad \eta : \text{learning rate}$$
$$\boldsymbol{\lambda} = \boldsymbol{\lambda}'$$

4. Iterate 1-3 until converge

Experiment - Multiple-Level Dur Models as PoE

Experimental conditions

- * Training data; 2,469 utterances
- * Development data; 137 utterances
 - Used to optimize weights in conventional method
 - Weights were optimized to minimize RMSE by grid search
 - Baseline & proposed method did not use development data
- * Almost the same training setup as Nitech-HTS 2005 [Zen;'06]
- * Test data; 137 utterances
- * State, phone, & syllable-level models were clustered individually
 - # of leaf nodes
 - * state; 607, phoneme; 1,364, syllable; 281

Experimental Results

Duration prediction results (RMSE in frame (rel imp))

Model	Phoneme	Syllable	Pause
Baseline (st)	5.08 (ref)	8.98 (ref)	35.0 (ref)
uPoE (st*ph)	4.62 (9.1%)	8.13 (9.5%)	31.8 (9.1%)
uPoE (st*ph*syl)	4.62 (9.1%)	8.11 (9.7%)	31.8 (9.1%)
PoE (st*ph)	4.60 (9.4%)	8.04 (10.5%)	31.9 (8.9%)
PoE (st*ph*syl)	4.57 (10.0%)	8.02 (10.7%)	31.9 (8.9%)

st; state only, st*ph; state & phoneme, st*ph*syl; state, phoneme, & syllable

uPoE; individually trained multiple-level duration models with optimized weights

PoE; jointly estimated multiple-level duration models

Experiment - Global Variance as PoE

Experimental conditions

- * Training data; 2,469 utterances
 - Training data was split into mini-batch (250 utterances)
 - 10 MCMC sampling at each contrastive divergence learning
 - * Hybrid Monte Carlo with 20 leap-frog steps
 - * Leap-frog size was adjusted adaptively
 - Learning rate was annealed at every 2,000 iterations
 - Momentum method was used to accelerate learning
 - Context-dependent logarithmic GV w/o silence was used
- * Test sentences; 70 sentences
 - Paired comparison test, # of subjects 7 (native English speaker)
 - 30 sentences per subject

Experimental Results

Paired comparison test result

Baseline	PoE	No preference
17.1	32.4	50.5

Baseline; conventional (not jointly estimated) GV

PoE; proposed (jointly estimated) GV

Difference was statistically significant at $p < 0.05$ level

Conclusions

Statistical parametric synthesis based on PoE

- **Combination of multiple-level AMs is formulated as PoE**
- **Jointly estimate multiple-level AMs as PoE**
 - * Linear feature function with Gaussian experts
 - Can be estimated in the same way as trajectory HMM
 - * Non-linear feature function and/or non-Gaussian experts
 - Contrastive divergence learning
- **Experiments**
 - * Jointly estimating multiple AMs as PoE improved performance

Future plans

- **Investigate other feature functions/experts**
- **Make contrastive divergence learning faster**

References

- [Hinton;'02] G. Hinton, "Training product of experts by minimizing contrastive divergence," *Neural Computation*, vol. 14, no. 8, pp. 1771-1800, 2002.
- [Williams;'05] C. Williams, "How to pretend that correlated variables are independent by using difference observations," *Neural Computation*, vol. 17, no. 1, pp. 1--6, 2005.
- [Zen;'07] H. Zen, et al., "Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic feature vector sequences," *Computer Speech & Language*, vol. 21, no. 1, pp.153-173, 2007.
- [Latorre;'08] J. Latorre & M. Akamine, "Multilevel parametric-base F0 model for speech synthesis," *Proc. Interspeech*, pp. 2274--2277, 2008.
- [Qian;'09] Y. Qian, et al., "Improved prosody generation by maximizing joint likelihood of state and longer units," *Proc. ICASSP*, pp. 1173--1176, 2009.
- [Wang;'08] C.-C. Wang, et al., "Multi-layer F0 modeling for HMM-based speech synthesis," *Proc. ISCSLP*, pp. 129--132, 2008.
- [Ling;'06] Z.-H. Ling, et al., "USTC system for Blizzard Challenge 2006 an improved HMM-based speech synthesis method," *Proc. Blizzard Challenge Workshop*, 2006.
- [Gao;'08] B.-H. Gao, et al., "Duration refinement by jointly optimizing state and longer unit likelihood," *Proc. Interspeech*, 2266--2269, 2008.
- [Toda;'07] T. Toda & K. Tokuda, "A speech parameter generation algorithm considering global variance for HMM-based speech synthesis," *IEICE Trans. Inf. Syst.*, vol. E90-D, no. 5, pp. 816--824, 2007.