T2R2 東京科学大学 リサーチリポジトリ Science Tokyo Research Repository

論文 / 著書情報 Article / Book Information

論題(和文)	GPR音声合成における深層構造の利用の検討				
Title(English)					
著者(和文)	都山知樹, 小林隆夫				
Authors(English)	Tomoki Koriyama, Takao Kobayashi				
出典(和文)	日本音響学会2018年春季研究発表会講演論文集, Vol. , No. , pp. 1507- 1508				
Citation(English)	, Vol. , No. , pp. 1507-1508				
 発行日 / Pub. date	2018, 3				

GPR 音声合成における深層構造の利用の検討*

○郡山知樹,小林隆夫(東工大)

1 はじめに

ガウス過程回帰 (GPR) に基づく音声合成は,カー ネル回帰のベイズ推定の枠組みで、コンテキスト同 士の相関関係を用いて音声パラメータを生成する手 法である [1]. GPR 音声合成においてカーネルの設計 は重要な課題であるが、例えばトライフォンのよう なカテゴリカルなコンテキストの類似度を表現する カーネルを、人手で設計することは困難である.近 年,ガウス過程において深層構造を用いて入力変数を 変換することによって、カーネルの構築を行う手法が 提案されている [2, 3]. これまでの報告 [4] で我々は, DNN を用いてコンテキストの特徴表現を抽出し、そ れを GPR の入力とする GP-DNN ハイブリッドモデ ル [2] の有効性を示した.本研究では, GP-DNN ハ イブリッドモデルだけでなく、ガウス過程自体を深層 構造に適用した深層ガウス過程 [3] を用いた音響モデ リング手法を提案し、GPR 音声合成における深層構 造の有用性を評価する.

2 GP-DNN ハイブリッドモデル

GPR 音声合成では、データ点i(i = 1...N)、次元 dの音響特徴量 y_i^d のノイズ成分を除いた潜在関数変 数 $f_i^d = f^d(\mathbf{x}_i)$ に対し、潜在関数 $f^d(\cdot)$ がガウス過程 $\mathcal{GP}(m^d(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'))$ に従うと仮定する.ただし、x は フレームレベルのコンテキスト、 $m^d(\mathbf{x})$ は平均関数、 $k(\mathbf{x}, \mathbf{x}')$ はコンテキスト同士の関係を表すカーネル関 数である.

GP-DNN ハイブリッドモデルに基づく音声合成 [4] では、Fig.1(b) に示すように、コンテキストをフィー ドフォワード型の DNN の入力とし、得られた低次元 特徴量 $\phi(\mathbf{x})$ をカーネル関数の入力とする. このとき、 確率的変分ガウス過程 (SVGP)[5] をガウス過程の近 似として用いると、以下の対数周辺尤度の変分下限 \mathcal{L} を最大化することで GP-DNN ハイブリッドモデル の学習が行われる.

$$\mathcal{L} = \sum_{i=1}^{N} \sum_{d=1}^{D} \left\{ \mathbb{E}_{q(f_i^d)} \left[\log p(y_i^d | f_i^d)) \right] - \frac{1}{N} \mathrm{KL}(q(\mathbf{u}^d) \| p(\mathbf{u}^d)) \right\}$$
(1)

$$q(f_i^d) = \mathcal{N}(f_i^d; \mu_i^d, \Sigma_i^d) \tag{2}$$

$$\mu_i^d = \mathbf{K}_{iM} \mathbf{K}_M^{-1} \mathbf{m}^d \tag{3}$$

 $\Sigma_i^d = \mathbf{K}_i - \mathbf{K}_{iM} \mathbf{K}_M^{-1} (\mathbf{K}_M - \mathbf{S}^d) \mathbf{K}_M^{-1} \mathbf{K}_{Mi}$ (4) ただし, $\mathbf{K}_{iM}, \mathbf{K}_i, \mathbf{K}_M$ は $k(\phi(\mathbf{x}), \phi(\mathbf{x}'))$ を要素に持 つグラム行列, $q(\mathbf{u}^d)$ は補助変数 \mathbf{u}^d の変分分布であ りガウス分布 $\mathcal{N}(\mathbf{u}^d; \mathbf{m}^d, \mathbf{S}^d)$ で表される. GP-DNN ハイブリッドモデルでは, 誤差逆伝播法により SVGP のモデルパラメータと DNN のパラメータを同時に学 習できる. また, 変分下限がデータ点ごとの和となっ ているため, 確率的勾配降下法に基づく最適化が可 能である. 音声合成時には, 式 (2) を用いて音声パラ メータの予測分布を推定し, これを用いて音声パラ



Fig. 1 DNN および提案法のモデル構造

メータ生成を行う.

3 深層ガウス過程

GP-DNN ハイブリッドモデルでは DNN の深層構 造によってコンテキストの特徴表現を抽出可能であ ることを想定しているが、重みの事前分布を考慮し ていないため DNN と同様に重みのパラメータが過 学習しやすい.これに対し、Fig.1(c) に示すような、 GPR を階層構造に接続した深層ガウス過程 (GP)[6] が提案されている.深層 GP は、層 *l*、次元 *d* の隠れ 関数 $f_i^{l,d}$ (*d* = 1...*D*_l) が *D*_{l-1} 次元の層 *l* – 1 の隠 れ関数変数 f_i^{l-1} を入力とするガウス過程に従って生 成されることを仮定したモデルである.深層 GP で は、すべての層で GPR によるベイズ推定を行うため、 GP-DNN に比べ過学習しにくいモデルといえる。

しかし深層 GP を計算量の問題から音声合成に直 接適用することが困難である.そこで本研究では, 大量のデータにも適用可能な手法として提案された DSVI(doubly stochastic variational inference) に基 づく深層 GP[3]を用いる.DSVIではSVGPと同様 に対数周辺尤度の変分下限を以下のようなサンプル 点ごとの値の和の形式で表す.

$$\mathcal{L} = \sum_{i=1}^{N} \left\{ \sum_{d=1}^{D_{L}} \mathbb{E}_{q(f_{i}^{L,d})} \left[\log p(y_{i}^{L,d} | f_{i}^{L,d}) \right] - \frac{1}{N} \sum_{l=1}^{L} \sum_{d=1}^{D_{l}} \mathrm{KL}(q(\mathbf{u}^{l,d}) \| p(\mathbf{u}^{l,d})) \right\}$$
(5)

ただし $q(\mathbf{u}^{l,d})$ は層l,次元dの補助変数 $\mathbf{u}^{l,d}$ の変分分 布であり、ガウス分布 $\mathcal{N}(\mathbf{u}^{l,d};\mathbf{m}^{l,d},\mathbf{S}^{l,d})$ で表される. 変分事後分布 $q(\mathbf{f}_{i}^{L})$ は

$$q(\mathbf{f}_i^L) = \int \prod_{l=1}^L q(\mathbf{f}_i^l | \mathbf{f}_i^{l-1}) d\mathbf{f}_i^l \tag{6}$$

となるが、この積分は一部のカーネル関数を除いて 解析的に計算不可能である。そこで DSVI では、下位 層のサンプル点 $\hat{\mathbf{f}}_{i}^{l-1}$ から上位層の隠れ変数 $\hat{\mathbf{f}}_{i}^{l}$ を条件 付き分布 $q(\mathbf{f}_{i}^{l}|\hat{\mathbf{f}}_{i}^{l-1})$ を用いて再帰的にサンプリングす る。S 個のサンプル点を用いると変分事後分布 $q(\mathbf{f}_{i}^{L})$ は、以下の混合ガウス分布で近似される。

$$q(\mathbf{f}_{i}^{L}) = \frac{1}{S} \sum_{s=1}^{S} q(\mathbf{f}_{i}^{L} | \hat{\mathbf{f}}_{i,s}^{L-1})$$
(7)

^{*}On the Use of Deep Model Architecture for GPR-based Speech Synthesis. by KORIYAMA, Tomoki, KOBAYASHI, Takao (Tokyo Institute of Technology)

Table 1 合成音声の原音声に対する音響特徴量歪 (MCEP: メルケプストラム距離 [dB], f_o: 対数 f_o の RMSE[cent], V/UV: 有声/無声誤り率 [%], BAP: 非周期性指標歪 [dB], DUR: 音素継続長歪 [ms])

手法	MCEP	$f_{\rm o}$	V/UV	BAP	DUR
DNN(ReLU)	4.80	175	4.70	3.10	16.4
DNN(tanh)	5.21	193	5.42	3.26	19.1
LSTM	4.71	183	4.65	3.08	18.0
GP-DNN	4.78	175	4.55	3.13	17.4
DeepGP	4.74	174	4.49	3.09	15.1

これを式 (5) に代入すると、変分下限はデータ点とサ ンプリングによって得られた点の二重の和の形式で 表される.したがって、SVGP と同様に DSVI に基 づく深層 GP では、確率的勾配降下法に基づくパラ メータ最適化が可能である.合成時においても学習時 と同様に、予測変分分布 $q(\mathbf{f}_i^L)$ を解析的に求められな い.本研究では、各層の予測平均 $\mathbb{E}[q(\mathbf{f}_i^l|\mathbf{\hat{f}}_i^{l-1})]$ を $\mathbf{\hat{f}}_i^l$ として予測分布 $q(\mathbf{f}_i^L)$ を近似する.

4 実験

4.1 実験条件

データベースには音声合成システム XIMERA[7] に 含まれる女性話者 F009 を使用した. 学習データには 1593 文 (約 119 分),評価データには 60 文 (約 4.1 分) の音声を用いた. 用いられた文には音素バランス文に 加え,旅行会話文,新聞読み上げ文が含まれている. サンプリングレート 16kHz の音声信号から,5ms 毎 に STRAIGHT を用いて f_o ,スペクトル包絡,非問 期性指標を抽出し,0-39 次のメルケプストラム,対 数 f_o ,5 次元の非周期性指標,およびそれらの Δ , Δ^2 と有声/無声フラグを音響特徴量として使用した.継 続長モデルでは音素継続長を音響特徴量とした.

音声単位ごとにコンテキストを抽出し, 音素, モー ラ, アクセント句, 呼気段落および発話のコンテキス トの次元を249, 88, 142, 82, 41 とした. GP-DNN ハイブリッドモデルおよび, 深層 GP の最上位層の カーネルには文献 [4] と同様に音声イベントごとの カーネルの和を求める加算構造のカーネルを用いる. SVGP の補助点数は1024 とした. またパラメータ最 適化手法には Adam を用いた.

GP-DNN ハイブリッドモデル (**GP-DNN**)の DNN の隠れ層は4層,256素子とし、出力層は32 次元とした.活性化関数には、勾配消失問題の起こり にくい手法として提案されている SELU[8]を使用し た.深層 GP(**DeepGP**)では層の数を6とし、最上 位層以外の層では補助点数は256,次元は32,カー ネルは RBF カーネルとした.

比較手法として DNN 音声合成 (**DNN**)[9] と双方 向 LSTM に基づく RNN 音声合成 (**LSTM**)[10] を用 いた. DNN 音声合成では隠れ層は 6 層, ノード数は 1024 とし, 活性化関数には ReLU および tanh を用い た. RNN 音声合成では下位 2 層のフィードフォワード 層と上位 2 層の双方向 LSTM 層からなるネットワー クを使用した. 活性化関数は tanh, 隠れ層の素子数 は 512 とした.

4.2 結果

客観評価結果として合成音声の原音声に対する歪 を Table 1 に示す.メルケプストラムや非周期性指標 は LSTM において最も歪が小さくなったが,対数 f_o,



Fig. 2 主観評価結果

有声/無声フラグ,音素継続長の歪は深層 GP が最も 小さいという結果を得た.また表から,DNN 音声合 成や GP-DNN ハイブリッドモデルに比べ深層 GP は すべての特徴量で歪が小さいことがわかる.

主観評価実験では MOS 試験による自然性の比較評価を行った. 被験者は 7 名で各被験者は評価データ 60 文の中からランダムに選ばれた 15 文を評価した. 合成音声の自然性を 5 段階で評価し,その MOS 値を 求めた. 結果を Fig.2 に示す. DNN 音声合成と他手 法を比較すると, a = 0.05 で有意に DNN の MOS の値が低かった. また, GP-DNN ハイブリッドモデ ルと深層 GP は,モデルにリカレント構造を有して いないにも関わらず,LSTM-RNN に基づく音声合成 と同程度のスコアを得た.

5 おわりに

本稿では GPR 音声合成に深層構造を導入した手法 として、GP-DNN ハイブリッドモデルと深層ガウス 過程に基づく手法を提案した.提案法は DNN 音声合成 より有意に自然性が高く、LSTM-RNN 音声合成と 同程度の自然性という主観評価結果となった.提案法 へのリカレント構造の導入や様々なデータでの評価 が今後の課題である.

謝 辞 本 研 究 は JSPS 科 研 費 JP15H02724, JP17K12711の助成を受けた.

参考文献

- [1] 郡山 他, "ガウス過程回帰に基づく音声合成システムの評価," 音講論 (秋), 3-1-3, pp. 235-236, 2015.
- [2] Bradshaw et al., "Adversarial examples, uncertainty, and transfer testing robustness in Gaussian process hybrid deep networks," arXiv:1707.02476, 2017.
- [3] Salimbeni et al., "Doubly stochastic variational inference for deep Gaussian processes," Proc. NIPS, pp.4591–4602, 2017.
- [4] 郡山 他, "GP-DNN ハイブリッドモデルに基づく統計的音声合成の検討,"信学技報, SP2017-67, 2018.
- Hensman et al., "Scalable variational Gaussian process classification," Proc. AISTATS, pp.1648–1656, 2015.
- [6] Damianou et al., "Deep Gaussian processes," Proc. AISTATS, pp.207–215, 2013.
- [7] 河井 他, "大規模コーパスを用いた音声合成システムXIMERA," 信学論 (D), 89(12), pp. 2688–2968, 2006.
- [8] Klambauer et al., "Self-normalizing neural networks," Proc. NIPS, pp.972–981, 2017.
- [9] Zen et al., "Statistical parametric speech synthesis using deep neural networks," Proc. ICASSP, pp.7962–7966, 2013.
- [10] Fan et al., "TTS synthesis with bidirectional LSTM based recurrent neural networks," Proc. IN-TERSPEECH, pp.1964–1968, 2014.