

波形接続型音声合成における歌唱の ピッチ変化予測 LSTM モデル

田中 瑞穂¹ 平田 圭二¹ 竹川 佳成¹

1. はじめに

近年、歌声合成システムを使って楽曲を制作するユーザーが増加している。このような歌声を生成するための音声合成ソフトは、表現力を高めるために歌声を調節することができる。しかし、人間歌唱のような物理的な制約がない合成歌唱は、ピッチや音量を無制限に変化させることができる。そのため、ソフト側でパラメータを設定することでポルタメントを一貫して制御している。一方で、制御が可能なパラメータが多いため、初心者は設定に時間がかかる。この問題を解決するため、機械学習によって人間歌唱を真似た合成歌唱生成する研究が注目されている。しかし、合成音声歌唱には、高速歌唱や急激な音程の変化、短いヴィブラートなど、人間では困難な合成音声特有の歌唱技術が多数存在している。これらの歌唱技術は、人間歌唱を用いる従来の手法では学習することができない。

そこで、本研究では、波形接続型合成音声を用いた合成音声ソフトで使用されているテキストデータを元にコーパスの作成を行い、未知の楽譜のパラメータを予測するモデルを提案する。この手法は、合成音声特有の歌唱技術のデータが含まれているため、人間歌唱では困難な歌唱技術にも対応が可能である。また、統計的合成音声ではなく波形接続型合成音声を用いた合成音声ソフトを選択した理由として、学習するパラメータが記号ベースであるため学習時間の短縮につながる事が挙げられる。なお、本研究では、抑揚を表現するパラメータに焦点を当てて学習する。

2. 合成音声ソフト UTAU

本研究では、合成音声ソフト UTAU [1] で使用される楽譜 UTAU Sequence Text(以下 UST とする) から抽出した音価、音高、歌詞、ポルタメントのデータをコーパスとする。

一般的にポルタメントは、音から音へ移動する際に、音程を滑らかに変化させる演奏及び歌唱技法を示す。一方、UTAU で用いるポルタメントは、1つの音の途中でも音程をずらすことができる。UTAU では、ポルタメント機能を用いることで図1のように点と点を繋ぐピッチ線を生成し、このピッチ線の形状に従って歌唱中の抑揚が変化する。なお、図1の赤い四角はピッチポイントと呼ばれる。UTAU のデフォルトのポルタメントは、2つのピッチポイントだけでスムーズに音を繋ぐものである。しかし、ピッチポイントの数を増やすことで、ピッチ線を曲げて細かく音程を変化させることが可能となる。

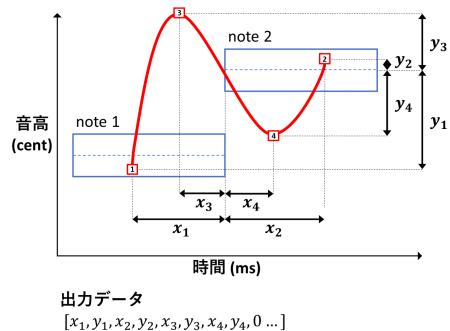


図 1 UTAU におけるポルタメントと出力例

3. パラメータ予測 LSTM モデル

本研究では、機械学習を用いてポルタメントに用いるパラメータを予測する。使用する機械学習モデルは、時系列を考慮しており音の前後関係が学習可能な LSTM [2] である。

入力データは、音価、音高、音のタイプ、前後の音高差を示す5つの One-hot ベクトルである。出力データは、ピッチポイントの個数を示す One-hot ベクトルとピッチポイントの座標を示すベクトルである。ピッチポイントの個数は、指定した音に対してヴィブラート及びポルタメントが

¹ 公立はこだて未来大学

出現するか判別するために用いる。

提案モデルは、大きく2つに分かれる。1つは、図2のように生成を判別するピッチポイントの個数のパラメータを他のパラメータとまとめて出力するモデルである。もう1つは、生成を判別するパラメータを別途学習するモデルである。これは、生成を判別するパラメータを一度学習してから入力に加えた方がまとめて学習するよりも精度が上がる考えたためである。

モデルの学習時は、1音ずつずらして楽譜全体を学習する。なお、本研究では、前の31音を関連付ける32次元のLSTMモデルと前の16音と後ろの16音を関連付ける33次元のLSTMモデルをそれぞれ用意している。これは、前の音のみを関連付けた時と前後両方の音を関連付けた時のどちらが精度が上がるか検証するためである。

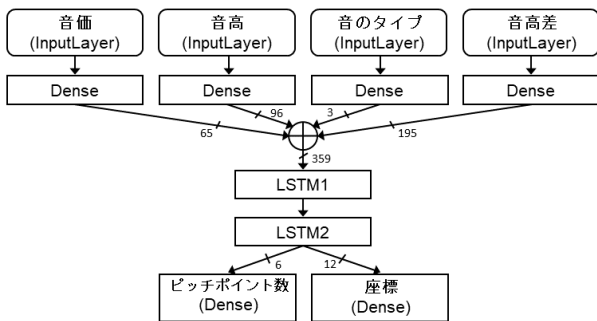


図2 予測 LSTM モデルのネットワーク

4. ポルタメント再現実験

本実験では、モデルがユーザの作成するパラメータと近い値を予測可能か検証する。モデルは、表1の4パターンを用いる。縦軸は関連付ける範囲、横軸は生成を判別するパラメータの扱いを示している。

	まとめて学習	別途学習
前31音	モデル1	モデル2
前16音及び後16音	モデル3	モデル4

まず、特定のユーザが作成したUSTデータを33曲分収集する。次に、集めたデータのうち30曲をコーパスとし、4章で解説したLSTMモデルで学習する。この際、エポック数は1000、バッチサイズは32とする。この学習モデルを使用し、正解データを元に作成された未知データの各パラメータを予測する。なお、本実験ではベースラインとして、UTAUの既存プラグイン機能であるAutoPitchwriterを使用している。

実験で評価する項目は、ポルタメントの出現箇所とポルタメントの形状の類似度の2点である。ポルタメントの出現箇所では、ピッチポイントが3個以上の抑揚をつけるポ

ルタメントが各音に付与されているか正解データと予測データで比較し正答率、適合率、再現率、F値を算出する。ポルタメントの形状の類似度では、正解データと予測データでポルタメントに何cent差があるか0.1msごとに算出し、その平均を取る。なお、ポルタメントが出現しない箇所及び出現する前後では、シフト値が0centであると仮定して比較する。

ポルタメントの出現箇所の実験結果は、表2のとおりである。ポルタメントの形状の類似度の実験結果は、表3のとおりである。類似度は、値が小さいほど正解データと誤差が少なく精度が高い結果となっている。

表2 ポルタメントの出現箇所

	Accuracy	Precision	Recall	F-measure
モデル1	65.0	61.0	99.0	75.3
モデル2	68.4	73.2	63.2	67.8
モデル3	73.2	82.0	64.2	71.7
モデル4	63.7	68.5	60.5	64.2
ベースライン	64.1	90.4	37.6	52.7

表3 ポルタメントの形状の類似度

モデル1	モデル2	モデル3	モデル4	ベースライン
10.0	10.6	10.2	12.8	12.1

5. 考察と今後の課題

表2、表3から、モデル1、3がモデル2、4に比べて出現箇所のF-measure及び形状の類似度が高くなっていった。このことより、ポルタメントの生成を判別するパラメータをまとめて学習するモデルの方が精度が高いと言える。また、モデル1、2とモデル3、4を比較した結果、モデル1、2の方が精度が高くなっていった。このことより、前の31音を関連付ける32次元のLSTMモデルの方が精度が高いと言える。

本稿では、合成音声ソフトに用いられるテキストデータを元にコーパスの作成を行い、未知の楽譜のパラメータを予測するモデルを提案した。今後は、入力に歌詞やフレーズなどのデータを与えることで精度の向上を目指す。

参考文献

- [1] 飴屋/菖蒲. 歌声合成ツール UTAU サポートページ. 入手先 (<http://utau2008.web.fc2.com/>) (参照 2021-12-21)
- [2] HOCHREITER, Sepp; SCHMIDHUBER, Jürgen. LSTM can solve hard long time lag problems. In: Advances in neural information processing systems. p. 473-479. 1997.