

# 敵対的DNN音声合成におけるダイバージェンスの影響の調査\*

◎齋藤 佑樹, 高道 慎之介, 猿渡 洋 (東大院・情報理工)

## 1 はじめに

統計的パラメトリック音声合成では、音響モデルから生成される合成音声パラメータ系列の過剰な平滑化により、音質の劣化が生じる。これまでに我々は、Deep Neural Network (DNN) 音声合成のための、Generative Adversarial Network (GAN) [1] の枠組みに基づく音響モデル学習 (敵対的DNN音声合成) [2] を提案し、その音質改善効果を確認している。GANの学習では、自然音声パラメータと合成音声パラメータ分布間距離を最小化するため、合成音声パラメータの分布の縮小を緩和することが可能である。

これまでに、Goodfellowら [1] によって提案されたGANを始めとして、数多くのGANが提案されている。しかし、これらのGANの有効性は主に画像生成においてのみ検証されており、同様の枠組みが敵対的DNN音声合成でも有効かは調査されていない。そこで、本稿では、各GANで最小化される距離規範の違いに着目し、それらが敵対的DNN音声合成の音質に与える影響を調査する。本稿では、音声信号処理との関連性が高い  $f$  ダイバージェンス最小化に基づくGAN ( $f$ -GAN) [3] と、画像生成において有効である Wasserstein GAN (W-GAN) [4], Least Squares GAN (LS-GAN) [5] を採用した敵対的DNN音声合成を比較する。実験的評価では、種々の距離規範の中で、Earth Mover 距離最小化に基づくW-GANが音質改善に最も有効であることを示す。

## 2 従来の枠組み

### 2.1 Minimum Generation Error (MGE) 学習による音響モデル学習

従来のMGE学習 [6] による音響モデル学習の損失関数  $L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}})$  は、自然音声のパラメータ系列  $\mathbf{y} = [\mathbf{y}_1^\top, \dots, \mathbf{y}_t^\top, \dots, \mathbf{y}_T^\top]^\top$  と、最尤パラメータ生成 [7] 後の合成音声のパラメータ系列  $\hat{\mathbf{y}} = [\hat{\mathbf{y}}_1^\top, \dots, \hat{\mathbf{y}}_t^\top, \dots, \hat{\mathbf{y}}_T^\top]^\top$  の間の二乗誤差として次式で与えられる。

$$\begin{aligned} L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}}) &= \frac{1}{T} (\hat{\mathbf{y}} - \mathbf{y})^\top (\hat{\mathbf{y}} - \mathbf{y}) \\ &= \frac{1}{T} (\mathbf{R}\hat{\mathbf{Y}} - \mathbf{y})^\top (\mathbf{R}\hat{\mathbf{Y}} - \mathbf{y}) \end{aligned} \quad (1)$$

ここで、 $t$  はフレームインデックス、 $T$  は総フレーム数、 $\mathbf{y}_t$  はフレーム  $t$  における音声パラメータ、 $\hat{\mathbf{Y}}$  は音響モデルの予測結果として得られる音声パラメータの静的・動的特徴量系列である。行列  $\mathbf{R}$  は、静的・動的特徴量の制約行列 [7] 及び学習データを用いて別途推定される分散・共分散行列  $\Sigma$  を用いて次式で計算される。

$$\mathbf{R} = (\mathbf{W}^\top \Sigma^{-1} \mathbf{W})^{-1} \mathbf{W}^\top \Sigma^{-1} \quad (2)$$

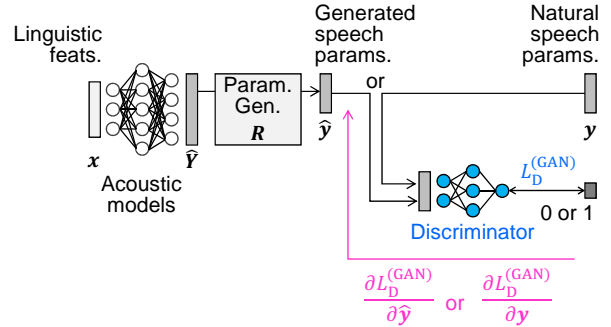


Fig. 1 識別モデルパラメータの更新時の計算フロー (Param. Gen. は最尤パラメータ生成 [7])

以降、この音声パラメータ生成を  $\hat{\mathbf{y}} = \mathbf{R}\hat{\mathbf{Y}} = \mathbf{G}(\mathbf{x}; \theta_{\text{G}})$  と定義する。 $\mathbf{x}$  は、入力テキストのコンテキスト系列である。音響モデルのモデルパラメータ  $\theta_{\text{G}}$  は、式 (1) の生成誤差の勾配  $\nabla_{\theta_{\text{G}}} L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}})$  を用いた backpropagation により更新される。勾配計算に必要な  $\nabla_{\hat{\mathbf{y}}} L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}})$  は  $\mathbf{R}^\top (\hat{\mathbf{y}} - \mathbf{y})/T$  として計算される [6]。

### 2.2 GAN [1] の学習

GANはDNNを用いた生成モデルを学習する手法であり、生成モデルと識別モデル  $D(\mathbf{y}; \theta_{\text{D}})$  の2つのDNNを学習する。ここで、 $\theta_{\text{D}}$  は識別モデルのモデルパラメータである。識別モデルの出力  $D(\mathbf{y})$  に sigmoid 関数を適用した値  $1/(1 + \exp(-D(\mathbf{y})))$  は、入力が真のデータである事後確率を表す。識別モデルは、真のデータに対して事後確率が1、生成モデルから生成されたデータに対して事後確率が0となるように学習される。一方で、生成モデルは識別モデルを詐称するように学習される。すなわち、生成データに対して識別モデルの事後確率が1となるように学習される。

GANの学習では、識別モデルと生成モデルが交互に更新される。まず、真のデータ  $\mathbf{y}$  と生成データ  $\hat{\mathbf{y}}$  を用いて、識別モデルの更新に用いる識別損失 (cross-entropy 関数)  $L_{\text{D}}^{(\text{GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  を計算する。

$$\begin{aligned} L_{\text{D}}^{(\text{GAN})}(\mathbf{y}, \hat{\mathbf{y}}) &= -\frac{1}{T} \sum_{t=1}^T \log \frac{1}{1 + \exp(-D(\mathbf{y}_t))} \\ &\quad - \frac{1}{T} \sum_{t=1}^T \log \frac{\exp(-D(\hat{\mathbf{y}}_t))}{1 + \exp(-D(\hat{\mathbf{y}}_t))} \end{aligned} \quad (3)$$

識別モデルのモデルパラメータ  $\theta_{\text{D}}$  は、式 (3) の識別損失の勾配  $\nabla_{\theta_{\text{D}}} L_{\text{D}}^{(\text{GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  を用いた backpropagation により更新される。Figure 1 に  $\theta_{\text{D}}$  の更新時の計算フローを示す。識別モデルの更新後、生成モデ

\* Experimental Investigation of Divergences in Adversarial DNN-Based Speech Synthesis, by SAITO, Yuki, TAKAMICHI, Shinnosuke, and SARUWATARI, Hiroshi (The University of Tokyo).

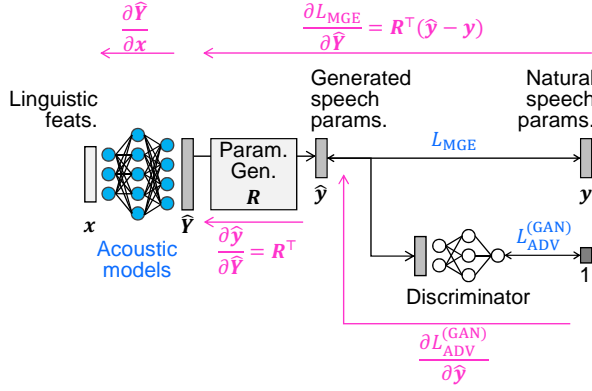


Fig. 2 音響モデルパラメータの更新時の計算フロー

の更新に用いる敵対損失  $L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}})$  を計算する.

$$L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T \log \frac{1}{1 + \exp(-D(\hat{\mathbf{y}}_t))} \quad (4)$$

生成モデルのモデルパラメータ  $\theta_G$  は, 式 (4) の敵対損失の勾配  $\nabla_{\theta_G} L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}})$  を用いた backpropagation により更新される. この枠組みにより, 真のデータの分布と生成されたデータの分布間の近似 Jensen-Shannon (JS) ダイバージェンスが最小化される [1].

### 2.3 Goodfellow らの GAN [1] を用いた敵対的 DNN 音声合成 [2]

Goodfellow らの GAN [1] を用いた敵対的 DNN 音声合成 [2] において, 音響モデル (生成モデル) の損失関数  $L_G(\mathbf{y}, \hat{\mathbf{y}})$  は次式で与えられる.

$$L_G(\mathbf{y}, \hat{\mathbf{y}}) = L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}}) + \omega_D \frac{E_{L_{\text{MGE}}}}{E_{L_{\text{ADV}}}} L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}}) \quad (5)$$

式 (5) の第二項  $L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}})$  は, 合成音声パラメータと自然音声パラメータの分布間距離 (近似 JS ダイバージェンス) を最小化する. 故に, 敵対的 DNN 音声合成は, パラメータの生成誤差  $L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}})$  を最小化するだけでなく, 合成音声パラメータの分布を自然音声のものに近づける.  $E_{L_{\text{MGE}}}$  と  $E_{L_{\text{ADV}}}$  はそれぞれ  $L_{\text{MGE}}(\mathbf{y}, \hat{\mathbf{y}})$  と  $L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}})$  の期待値であり, 2つの損失のスケールを調整する役割を持つ.  $\omega_D$  は, 敵対損失の影響を調整するハイパーパラメータである. 音響モデルのモデルパラメータ  $\theta_G$  は, 式 (5) の勾配  $\nabla_{\theta_G} L_{\text{ADV}}^{(\text{GAN})}(\hat{\mathbf{y}})$  を用いた backpropagation により更新される. Figure 2 に  $\theta_G$  の更新時の計算フローを示す. GAN の学習と同様に, 音響モデルと識別モデルは交互に更新され, 一方の更新中に他方のモデルパラメータは固定される.

## 3 敵対的 DNN 音声合成に用いる種々の距離規範

GAN の枠組みは真のデータと生成データの分布間の距離規範最小化とみなされる. 2.2 節で述べたように, Goodfellow らの GAN [1] は近似 JS ダイバージェンスを最小化する. この分布間距離最小化の観点に基づき,

様々な GAN を敵対的 DNN 音声合成に導入する. 本稿では, 非負値行列因子分解 [12, 13] における損失関数としてしばしば採用される Kullback-Leibler (KL) ダイバージェンスを含む  $f$  ダイバージェンスの最小化に基づく GAN ( $f$ -GAN) [3] と, 画像生成において有効である Wasserstein GAN (W-GAN) [4] 及び Least Squares GAN (LS-GAN) [5] を導入する. 以降で導入する識別モデルの識別損失  $L_D^{(*-\text{GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  と, 音響モデルの敵対損失  $L_{\text{ADV}}^{(*-\text{GAN})}(\hat{\mathbf{y}})$  は, それぞれ式 (3) と式 (4) の代わりに用いられる.

### 3.1 $f$ -GAN [3]

$f$ -GAN は, Goodfellow らの GAN [1] を包含する統一的な枠組みを記述する. 分布間の距離は, KL ダイバージェンスや JS ダイバージェンスを含む  $f$ -ダイバージェンス  $\mathcal{D}_f(\mathbf{y} \parallel \hat{\mathbf{y}})$  により次式で定義される.

$$\mathcal{D}_f(\mathbf{y} \parallel \hat{\mathbf{y}}) = \int q(\hat{\mathbf{y}}) f\left(\frac{p(\mathbf{y})}{q(\hat{\mathbf{y}})}\right) d\mathbf{y} \quad (6)$$

ここで,  $p(\cdot)$  と  $q(\cdot)$  はそれぞれ  $\mathbf{y}$  と  $\hat{\mathbf{y}}$  の確率密度である.  $f(\cdot)$  は  $f(1) = 0$  を満たす凸関数である.  $f(\cdot)$  の選択により様々なダイバージェンスを表すことが可能だが, 本稿では音声信号処理との関連性が高いものを採用する.

#### 3.1.1 KL-GAN

KL ダイバージェンスは,  $f(r) = r \log r$  と定義することで次式で与えられる.

$$\mathcal{D}_{\text{KL}}(\mathbf{y} \parallel \hat{\mathbf{y}}) = \int p(\mathbf{y}) \log \frac{p(\mathbf{y})}{q(\hat{\mathbf{y}})} d\mathbf{y} \quad (7)$$

識別モデルの識別損失  $L_D^{(\text{KL-GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  は

$$L_D^{(\text{KL-GAN})}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T D(\mathbf{y}_t) + \frac{1}{T} \sum_{t=1}^T \exp(D(\hat{\mathbf{y}}_t) - 1) \quad (8)$$

であり, 音響モデルの敵対損失  $L_{\text{ADV}}^{(\text{KL-GAN})}(\hat{\mathbf{y}})$  は

$$L_{\text{ADV}}^{(\text{KL-GAN})}(\hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T D(\hat{\mathbf{y}}_t) \quad (9)$$

として与えられる.

#### 3.1.2 Reversed KL (RKL)-GAN

KL ダイバージェンスは非対称な距離規範であるため, その双対な距離規範である reversed KL (RKL) ダイバージェンス  $\mathcal{D}_{\text{RKL}}(\mathbf{y} \parallel \hat{\mathbf{y}})$  は  $\mathcal{D}_{\text{KL}}(\mathbf{y} \parallel \hat{\mathbf{y}})$  と異なる. RKL ダイバージェンスは,  $f(r) = -\log r$  と定義することで次式で与えられる.

$$\mathcal{D}_{\text{RKL}}(\mathbf{y} \parallel \hat{\mathbf{y}}) = \int q(\hat{\mathbf{y}}) \log \frac{q(\hat{\mathbf{y}})}{p(\mathbf{y})} d\mathbf{y} = \mathcal{D}_{\text{KL}}(\hat{\mathbf{y}} \parallel \mathbf{y}) \quad (10)$$

識別モデルの識別損失  $L_D^{(\text{RKL-GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  は

$$L_D^{(\text{RKL-GAN})}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{T} \sum_{t=1}^T \exp(-D(\mathbf{y}_t)) + \frac{1}{T} \sum_{t=1}^T (-1 + D(\hat{\mathbf{y}}_t)) \quad (11)$$

であり、音響モデルの敵対損失  $L_{\text{ADV}}^{(\text{RKL-GAN})}(\hat{\mathbf{y}})$  は

$$L_{\text{ADV}}^{(\text{RKL-GAN})}(\hat{\mathbf{y}}) = \frac{1}{T} \sum_{t=1}^T \exp(-D(\hat{\mathbf{y}}_t)) \quad (12)$$

として与えられる。

### 3.1.3 JS-GAN

近似なしの JS ダイバージェンスの最小化は、 $f$ -GAN の枠組みで記述可能である。JS ダイバージェンスは、 $f(r) = -(r+1) \log \frac{r+1}{2} + r \log r$  と定義することで次式で与えられる。

$$\mathcal{D}_{\text{JS}}(\mathbf{y} \parallel \hat{\mathbf{y}}) = \frac{1}{2} \int p(\mathbf{y}) \log \frac{2p(\mathbf{y})}{p(\mathbf{y}) + q(\hat{\mathbf{y}})} d\mathbf{y} + \frac{1}{2} \int q(\hat{\mathbf{y}}) \log \frac{2q(\hat{\mathbf{y}})}{p(\mathbf{y}) + q(\hat{\mathbf{y}})} d\mathbf{y} \quad (13)$$

識別モデルの識別損失  $L_D^{(\text{JS-GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  は

$$L_D^{(\text{JS-GAN})}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T \log \frac{2}{1 + \exp(-D(\mathbf{y}_t))} - \frac{1}{T} \sum_{t=1}^T \log \frac{2 \exp(-D(\hat{\mathbf{y}}_t))}{1 + \exp(-D(\hat{\mathbf{y}}_t))} \quad (14)$$

であり、音響モデルの敵対損失  $L_{\text{ADV}}^{(\text{JS-GAN})}(\hat{\mathbf{y}})$  は

$$L_{\text{ADV}}^{(\text{JS-GAN})}(\hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T \log \frac{2}{1 + \exp(-D(\hat{\mathbf{y}}_t))} \quad (15)$$

として与えられる。Goodfellow らの GAN [1] は、 $2\mathcal{D}_{\text{JS}}(\mathbf{y} \parallel \hat{\mathbf{y}}) - \log(4)$  を最小化する。

### 3.2 Wasserstein GAN (W-GAN) [4]

Goodfellow らの GAN [1] の学習の不安定さを改善させるために、W-GAN が提案されている。W-GAN の学習では、次式に示す Earth-Mover 距離 (Wasserstein-1) を最小化する。

$$\mathcal{D}_{\text{EM}}(\mathbf{y}, \hat{\mathbf{y}}) = \inf_{\gamma} \mathbb{E}_{(\mathbf{y}, \hat{\mathbf{y}}) \sim \gamma} [\|\mathbf{y} - \hat{\mathbf{y}}\|] \quad (16)$$

ここで、 $\gamma(\mathbf{y}, \hat{\mathbf{y}})$  は周辺分布がそれぞれ  $\mathbf{y}$  と  $\hat{\mathbf{y}}$  の分布となるような結合分布である。Kantorovich-Rubinstein の定理 [14] を用いると、識別モデルの識別損失  $L_D^{(\text{W-GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  は

$$L_D^{(\text{W-GAN})}(\mathbf{y}, \hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T D(\mathbf{y}_t) + \frac{1}{T} \sum_{t=1}^T D(\hat{\mathbf{y}}_t) \quad (17)$$

であり、音響モデルの敵対損失  $L_{\text{ADV}}^{(\text{W-GAN})}(\hat{\mathbf{y}})$  は

$$L_{\text{ADV}}^{(\text{W-GAN})}(\hat{\mathbf{y}}) = -\frac{1}{T} \sum_{t=1}^T D(\hat{\mathbf{y}}_t) \quad (18)$$

として与えられる。W-GAN では、識別モデルが  $K$ -Lipschits 関数である必要がある。そのため、識別モデルの更新後に、重みパラメータの値を  $[-0.01, 0.01]$  のような一定区間内に収める必要がある。

### 3.3 Least Squares GAN (LS-GAN) [5]

Goodfellow らの GAN [1] では、sigmoid cross-entropy を用いるため、勾配消失が起きやすい。これを防ぐために、LS-GAN [5] では、目的関数を二乗誤差として定式化する。識別モデルの識別損失  $L_D^{(\text{LS-GAN})}(\mathbf{y}, \hat{\mathbf{y}})$  は

$$L_D^{(\text{LS-GAN})}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2T} \sum_{t=1}^T (D(\mathbf{y}_t) - b)^2 + \frac{1}{2T} \sum_{t=1}^T (D(\hat{\mathbf{y}}_t) - a)^2 \quad (19)$$

であり、音響モデルの敵対損失  $L_{\text{ADV}}^{(\text{LS-GAN})}(\hat{\mathbf{y}})$  は

$$L_{\text{ADV}}^{(\text{LS-GAN})}(\hat{\mathbf{y}}) = \frac{1}{2T} \sum_{t=1}^T (D(\hat{\mathbf{y}}_t) - c)^2 \quad (20)$$

として与えられる。ここで、 $a, b, c$  はそれぞれ、識別モデルに合成音声を合成音声と、自然音声を自然音声と、合成音声を自然音声として識別させるためのラベルである。これらのラベルが  $b - c = 1$  及び  $b - a = 2$  を満たすとき、LS-GAN は  $p(\mathbf{y}) + q(\hat{\mathbf{y}})$  と  $2q(\hat{\mathbf{y}})$  の間の Pearson  $\chi^2$  ダイバージェンスを最小化する。予備実験により、このダイバージェンスは合成音声の品質を劣化させることが確認されたため、本稿では [5] の式 (9) で提案されている  $a = 0, b = 1, c = 1$  を用いる。

## 4 実験的評価

### 4.1 実験条件

実験的に用いるデータとして、ATR 音素バランス 503 文 [15] を利用し、A-I セット 450 文を学習に、J セット 53 文を評価に用いる。学習データのサンプリング周波数は 16 kHz、フレームシフトは 5 ms とする。スペクトル特徴量として STRAIGHT 分析 [16] による 0 次から 24 次のメルケプストラム係数、音源特徴量として  $F_0$ 、5 周波数帯域における平均非周期成分 [17] を用いる。スペクトル特徴量に対する前処理として、50 Hz のカットオフ変調周波数による trajectory smoothing [18] を利用する。継続長は、自然音声のものを用いる。なお、0 次のメルケプストラム係数については、敵対的 DNN 音声合成による品質の劣化が確認されたため、MGE 学習後のものを用いる。コンテキストラベルは、音素、モーラ位置、アクセント型、音素内フレーム位置などから成る 442 次元ベクトルである。DNN 学習時には、音声特徴量及び実数値をとるコンテキストラベル特徴量を平均 0、

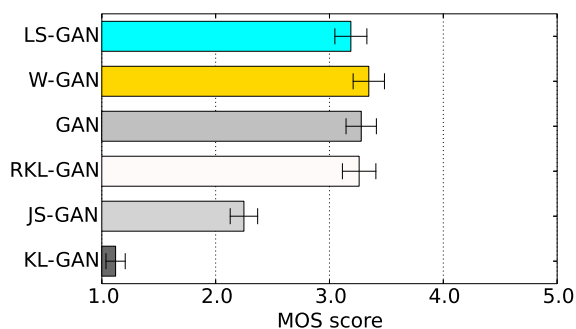


Fig. 3 音質に関する主観評価結果（エラーバーは95%信頼区間）

分散1に正規化する。音声合成の音響モデルと識別モデルのためのDNNはFeed-Forward型とする。音声合成の音響モデルの隠れ層数は3、隠れ層の素子数は512、隠れ層及び出力層の活性化関数はそれぞれRectified Linear Unit (ReLU) 及び線形関数である。識別モデルの隠れ層数は3、隠れ層の素子数は256、隠れ層の活性化関数はReLUである。識別モデルの出力層の活性化関数は、各GANにより異なる。音響モデルはスペクトル特徴量(25次元)、連続対数 $F_0$ (1次元)、非周期成分(5次元)とそれらの動的特徴量( $\Delta$ ,  $\Delta\Delta$ ), U/V(1次元)を結合させた94次元のベクトルをフレーム毎に予測し、識別モデルはスペクトル特徴量と連続対数 $F_0$ の静的特徴量を結合させた26次元のベクトルを用いて自然音声と合成音声をフレーム毎に識別する。最適化手法として学習率0.001のAdaGrad [19]を用いる。

まず、音響モデルの初期化として、反復回数25回のMGE学習を行う。次に、識別モデルの初期化として、自然音声特徴量とMGE学習後の合成音声特徴量を識別するような反復回数5回の学習を行う。その後、初期化された音響モデルと識別モデルを用いて、反復回数25回の敵対的DNN音声合成の学習を行う。

本稿では、以下のGANを採用した敵対的DNN音声合成を比較する。ただし、ハイパーパラメータ $\omega_D$ は、全てのGANで通常設定の1.0とする。

**GAN:** 式 (3), (4)

**KL-GAN:** 式 (8), (9)

**RKL-GAN:** 式 (11), (12)

**JS-GAN:** 式 (14), (15)

**W-GAN:** 式 (17), (18)

**LS-GAN:** 式 (19), (20)

主観評価として、我々のクラウドソーシングによる評価システムを用いて、55名の被験者による5段階のMOSテストを実施する。評価指標は、合成音声の音質である。

#### 4.2 評価結果

Figure 3に評価結果を示す。評価結果より、KL-GANとJS-GANを除く種々の距離規範が敵対的

DNN音声合成に有効であることが確認できる。興味深い点として、(1) KLダイバージェンスを最小化するKL-GANは有効でないが、その双対であるRKL-GANは有効であり、高いスコアを獲得していること、(2) JSダイバージェンスを最小化するJS-GANはあまり有効ではないが、その近似を最小化するGANは有効であることが挙げられる。最も高いスコアを獲得したGANはW-GANであり、LS-GAN、JS-GAN、そしてKL-GANと比較して有意に高い音質であることが分かった。この要因として、特に無声音及び無音区間における音質の改善を確認した。

## 5 おわりに

本稿では、種々の距離規範を最小化するGANを敵対的DNN音声合成に導入し、実験的評価によりその影響を調査した。評価結果より、Earth-Mover距離最小化に基づくWasserstein GANが、合成音声の品質改善に最も有効であること示した。今後は、スペクトログラムの生成 [20] において有効な距離規範及びDNNアーキテクチャについて調査する。

**謝辞** 本研究は、総合科学技術・イノベーション会議による革新的研究開発推進プログラム (ImPACT)、セコム科学技術振興財団、及びJSPS科研費16H06681の支援を受けた。

## 参考文献

- [1] Goodfellow et al., *Proc. NIPS*, pp. 2672–2680, 2014.
- [2] Saito et al., *Proc. ICASSP*, pp. 4900–4904, 2017.
- [3] Nowozin et al., *Proc. NIPS*, pp. 271–279, 2016.
- [4] Arjovsky et al., *arXiv:1701.07875*, 2017.
- [5] Mao et al., *arXiv:1611.04076*, 2017.
- [6] Wu et al., *IEEE/ACM Trans. on ASLP*, Vol. 24, No. 7, pp. 1255–1265, 2016.
- [7] Tokuda et al., *Proc. ICASSP*, pp. 1315–1318, 2000.
- [8] Huang et al., *Proc. INTERSPEECH*, pp. 2464–2468, 2015.
- [9] Reed et al., *Proc ICML*, pp. 1060–1069, 2016.
- [10] Wu et al., *IEEE/ACM Trans. on ASLP*, Vol. 24, No. 4, pp. 768–783, 2016.
- [11] Chen et al., *Proc. INTERSPEECH*, pp. 2097–2101, 2015.
- [12] Lee et al., *Proc. NIPS*, pp. 556–562, 2000.
- [13] Kompass, *Neural Computation*, Vol. 19, No. 3, pp. 780–891, 2007.
- [14] Vilani, *Optimal Transport: Old and New*, Springer, 2009.
- [15] 阿部 他, ATRテクニカルレポート, TR-I-0166, 1990.
- [16] Kawahara et al., *Speech Commun.*, Vol. 27, No. 3–4, pp. 187–207, 1999.
- [17] Ohtani et al., *Proc. INTERSPEECH*, pp. 2266–2269, 2006.
- [18] Takamichi et al., *Proc. Blizzard Challenge Workshop*, 2015.
- [19] Duchi et al., *JMLR*, Vol. 12, pp. 2121–2159, 2011.
- [20] Wang et al., *arXiv:1703.10135*, 2017.