

# GANを用いたマルチモーダル情報の生成

藤元陸<sup>†</sup> 堀井隆斗<sup>†</sup> 日永田智絵<sup>†</sup> 宮澤和貴<sup>†</sup> 長井隆行<sup>†</sup>  
電気通信大学<sup>†</sup>

## 1 はじめに

深層学習技術の発展に伴い、画像や音などの高次元データを違和感なく生成可能なモデルが提案されている<sup>2)</sup>。一方で、動画(画像と音声)や絵本(画像と文章)のような複数のモダリティで構成されるデータの生成に関してははまだ十分でない。マルチモーダルデータの生成は単一モダリティデータの生成と比べて困難な課題である。なぜなら、モダリティが増えることによりデータが複雑になり、また、モダリティ間の関係性を考慮して生成する必要があるからである。

本稿ではマルチモーダルデータの生成モデルを提案する。提案モデルは敵対的生成ネットワーク(GAN)<sup>1)</sup>を拡張して、複数データを生成できるようにしたモデルである。そして、このモデルで学習を行い、実際にマルチモーダルデータを生成できるか確認する。また、マルチモーダルデータはモダリティ間に関係性があるが、生成モデルが関係性を考慮してデータを生成しているかどうかを検証する。

## 2 Multimodal Generative Adversarial Networks(MGAN)

### 2.1 MGANの定式化

マルチモーダルデータを同時に生成するためのネットワーク、Multimodal Generative Adversarial Networks(MGAN)を提案する。Fig. 1は、生成データのモダリティが二種類の場合のMGANのネットワーク図である。MGANは複数のモダリティのデータを生成するために、複数のGANで構成される。通常のGANはノイズと生成データは1対1で対応するが、MGANではノイズは複数モダリティのデータと対応する。また、それぞれのGANのGeneratorの入力層とDiscriminatorの出力層は共有する。これにより、高次の特徴が共有される。

GAN1の生成データを $g_1(z)$ 、GAN2の生成データを $g_2(z)$ とすると、Discriminatorの出力は $d(g_1(z), g_2(z))$ と表される。このとき、MGANの学習は以下のミニマックスゲームで与えられる。

$$\max_{g_1, g_2} \min_d V(d, g_1, g_2) = E_{x \sim p_x}[-\log d(x)] + E_{z \sim p_z}[-\log d(1 - d(g_1(z), g_2(z)))] \quad (1)$$

ただし、 $x$ はデータセットの分布 $p_x$ からサンプルしたデータ $x$ を表す。 $z$ は適当なノイズ $z$ の分布 $p_z$ からサンプルしたノイズ $z$ を表す。

### 2.2 モダリティ間のバランス調整

MGANでは複数のモダリティを生成するが、データによって各モダリティの生成しやすさが異なる。その

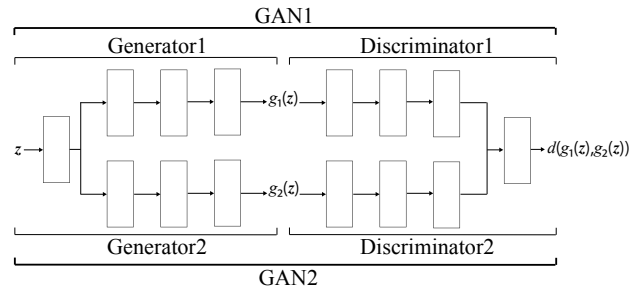


Fig. 1: Multimodal Generative Adversarial Networks

ため、あるモダリティは生成できるが、他のモダリティは生成できないということが発生する。そこで、モダリティごとの生成優先度の調整をFeature matching<sup>3)</sup>を用いて行う。Feature matchingを利用することで、ネットワーク構造を変化させずにモダリティ間の損失のバランスを調整できる。

Feature matchingはGANの生成を安定化させるテクニックの一つであり、Discriminatorに真のデータ $x$ を入力したときの中間層の出力 $f(x)$ と、 $g(z)$ を入力したときの中間層の出力 $f(g(z))$ のそれぞれの期待値が一致するようにGeneratorを訓練する。通常のGANでは以下の項をGeneratorの損失に加える。

$$\|E_{x \sim p_x} f(x) - E_{z \sim p_z} f(g(z))\|_2^2 \quad (2)$$

MGANでは各モダリティに対しFeature matchingを適用する。Discriminator1の中間層の出力を $f_1(\cdot)$ 、Discriminator2の中間層の出力を $f_2(\cdot)$ とすると以下の項をGeneratorの損失に加える。

$$a \|E_{x \sim p_x} f_1(x) - E_{z \sim p_z} f_1(g(z))\|_2^2 + b \|E_{x \sim p_x} f_2(x) - E_{z \sim p_z} f_2(g(z))\|_2^2 \quad (3)$$

$a, b$ はパラメータであり、式(1)に対するFeature matchingの重みと各モダリティの重みを調整する。 $a$ の重みを大きくするとモダリティ1に関する損失が大きくなり、モダリティ1の方が優先的にパラメータ更新される。逆に、 $b$ の重みを大きくするとモダリティ2のパラメータ更新が優先される。 $a$ と $b$ の重みはデータやネットワーク構造によって最適な値が変わるため、設計者が生成結果を考慮しつつ調整する必要がある。

## 3 実験

MGANがモダリティ間の対応関係を考慮してマルチモーダルデータを生成できるかの検証を行う。

### 3.1 データセット

モダリティ間の対応関係を検証するためにFig. 2(a)のロボットの腕をシミュレーション上で動かし、視覚情報と関節角情報を取得した。視覚情報は頭部のカメラから画像を取得し、画像サイズを $80 \times 80$ pixelに変

<sup>†</sup>Generation of multimodal information using GAN]

<sup>†</sup>The University of Electro-Communications



Fig. 3: 生成画像と関節角の再構成画像.

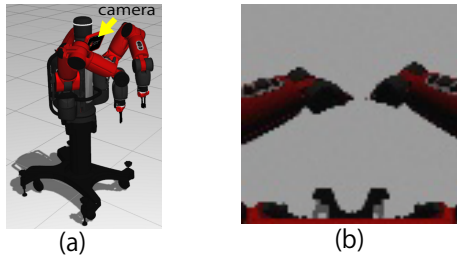


Fig. 2: 実験で用いたロボットと視覚画像

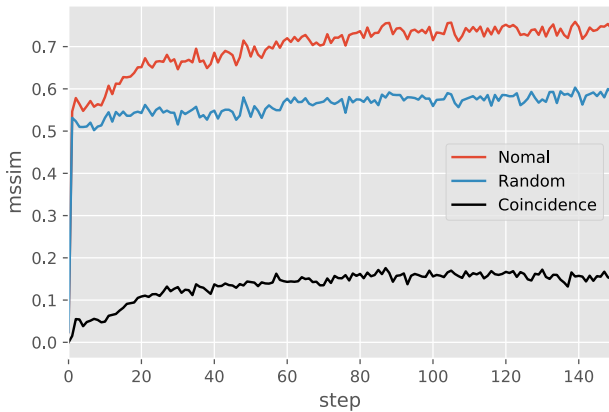


Fig. 4: SSIM result

換した。Fig. 2(b) は変換後の画像の一つである。関節角情報は片腕につき5関節を動かし、両腕で計10関節のデータを視覚情報とほぼ同時に取得した。

### 3.2 学習結果

MGANの学習結果をFig. 3に示す。上段の画像はMGANが生成した視覚情報、下段はMGANが生成した関節角情報を用いてシミュレーション上のロボットを動かして再現した画像である。視覚情報は綺麗な画像を生成できた。また、再現画像と生成した視覚情報を比較すると全体的な見た目や腕の位置が近いことがわかる。これは、視覚情報と関節角情報の対応関係を正しく学習できたことを示している。

### 3.3 SSIMによるモデルの評価

MGANがモダリティ間の関係性を考慮してデータを生成しているかをSSIM<sup>4)</sup>を用いて検証する。SSIMは画像間の類似度の手法である。生成画像と関節角情報による再現画像を学習中の各stepで10組生成し、stepごとに算出したSSIMの平均値mssimがFig. 4の赤線(Nomal)である。赤線は生成画像と再現画像がどれだ

け似ているかを表している。しかし、この値はあくまで画像の一致度であり、モダリティ間の一致度を表しているわけではない。モダリティ間の関係性を考慮して生成しているかを評価するためには、背景やロボットの動かない部位の影響を除く必要がある。そこで、各ステップで生成画像と再現画像のランダムな組み合わせを1000組作り、SSIMを算出した(Fig. 4の青線)。この値はモダリティ間の関係性の影響を除いた評価値だと考えられる。この二つの評価値の差によってモダリティ間の関係性を評価した(Fig. 4の黒線)。80step程で0.17程度に収束しており、MGANがモダリティ間の関係性を考慮して生成していることを示している。

## 4 結論

本稿では、マルチモーダルデータを生成するネットワークMGANを提案した。実際にシミュレータ上でロボットの腕を動かしたときの画像と関節角の情報を用いて学習し、マルチモーダルデータの生成を行った。生成した関節角情報をもとにシミュレータで再現した画像と生成画像を比較すると類似していた。また、SSIMをもちいてモダリティ間の共起性を算出したところ、学習と共に増加していた。このことから、MGANはモダリティ間の共起性を考慮してマルチモーダルデータの生成ができるネットワークであるといえる。今後の課題として、他の様々なモダリティへの適用、潜在空間との関連性の解析、時系列データへの対応などが挙げられる。

## 謝辞

本研究の一部は、JST CREST (JPMJCR15E3)の支援を受けて実施したものである。

## 参考文献

- 1) Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- 2) Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.
- 3) Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. In *Advances in Neural Information Processing Systems*, pp. 2234–2242, 2016.
- 4) Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, Vol. 13, No. 4, pp. 600–612, 2004.