

Generative Adversarial Networks

小泉 孝弥

1 Introduction

皆さんは、GAN というものを知っているだろうか。GAN^{*1}とは敵対的生成ネットワーク (Generative Adversarial Networks) のことであり、2014 年に登場した人工知能の技術である。「敵対的」という単語がついている通り、2 つのニューラルネットワークという関数を競わせることで、学習させていく手法である。今回の方程では、まず、機械学習についての説明をする。その後、GAN の数学的な定式化の説明をし、最後に GAN の学習を安定させる方法について述べる。

2 Fundamental concepts of Machine Learning

この節では機械学習の数学的な定式化を行う。 $(\Omega, \mathcal{F}, \mathbb{P})$ を完備な確率空間とする。 d, m を自然数とする。

Definition 2.1. (仮説空間, 仮説)

$\mathcal{X} := \mathbb{R}^d$ から $\mathcal{Y} := \mathbb{R}^m$ へのなんらかの条件を満たす写像の集まりのことを仮説空間といい \mathcal{H} と表記する。すなわち、

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f \text{ が満たす条件} \}$$

である。(条件の具体例は後に述べる)。仮説空間 \mathcal{H} の元のことを仮説 (またはモデル) と呼ぶ。また、この時の \mathcal{X} を入力空間、 \mathcal{Y} を出力空間と呼ぶ。

Definition 2.2. (データ)

\mathcal{X} と \mathcal{Y} の直積集合 $\mathcal{X} \times \mathcal{Y}$ の有限部分集合のことを教師ありデータという。また、 \mathcal{X} の有限部分集合のことを教師なしデータと呼ぶ。

Remark 2.3. (データについて)

データ D は確率分布 $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ にしてがって生成されているものとする。すなわち、 $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ に従う確率変数列 $\{(X_n, Y_n)\}_{n=1}^N$ の観測値のことをデータと呼ぶ。

^{*1} 「ぎゃん」と読む人と、「がん」と読む人がいるが、ここでは「ぎゃん」とすることにする。

Definition 2.4. (予測損失)

\mathcal{H} を仮説空間とし, (X, Y) をデータの分布 $\mathbb{P}_{\mathcal{X} \times \mathcal{Y}}$ に従う確率変数とする. 以下で定義される $\ell: \mathcal{H} \rightarrow \mathbb{R}$ を予測損失と呼ぶ.

$$\ell(f) = \mathbb{E}[\|f(X) - Y\|].$$

予測損失が最小となるような $f \in \mathcal{H}$ を求めることが機械学習の目標である. しかし, 一般にデータの分布は未知であるためこの式を解くことができない. そこで, 損失関数というものを定義し, それを最適化することを考える.

Definition 2.5. (損失関数)

\mathcal{H} を仮説空間とする. \mathcal{H} から \mathbb{R} への写像 $\mathcal{L}_D: \mathcal{H} \rightarrow \mathbb{R}$ を損失関数 (Loss function) と呼ぶ.

Definition 2.6. (機械学習空間)

D をデータ, \mathcal{H} を仮説空間, \mathcal{L}_D を損失関数とする. この時 5 つ組 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を機械学習空間 (Machine Learning space) という.

Definition 2.7. (学習, 最適仮説)

\mathcal{H} を仮説空間, $\mathcal{L}_D: \mathcal{H} \rightarrow \mathbb{R}$ を損失関数とする. 損失関数が最大または, 最小となるような $f^* \in \mathcal{H}$ を求めること^{*2}を学習といい, $f^* \in \mathcal{H}$ を最適仮説と呼ぶ.

Definition 2.8. (機械学習)

機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上での学習を機械学習という. 特に, D が教師ありデータの時, 教師あり機械学習と呼び, 教師なしデータの時, 教師なし機械学習と呼ぶ.

Remark 2.9. (訓練データとテストデータ)

機械学習で学習を行う際は, データ D を全て使って学習させるのではなく, データを訓練データ D_{Train} とテストデータ D_{Test} に分ける. 訓練データはモデルの学習に用いて, テストデータはモデルの性能の確認のために用いる.

3 教師あり機械学習の具体例

前節では, 機械学習の抽象的な枠組みを紹介したが, この説では機械学習空間の具体例を述べる. $D = \{(x_i, y_i)\}_{i=1}^N \subset \mathcal{X} \times \mathcal{Y}$ を教師ありデータとする.

^{*2} 機械学習においては厳密に損失関数が最大・最小となる関数が学習されるとは限らないが, そのような関数を求めることも学習の定義に含めるものとする.

3.1 単回帰分析・重回帰分析

最初に機械学習の最も基本的なモデルである単回帰分析を紹介する．その後その多変数拡張である重回帰分析を紹介する．

Example 3.1. (単回帰分析)

機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を以下のように定義する．

$$\mathcal{X} = \mathbb{R}, \mathcal{Y} = \mathbb{R},$$

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f(x) = wx, w \in \mathbb{R}\},$$

$$\mathcal{L}_D(f) = \sum_{i=1}^N |f(x_i) - y_i|^2.$$

この機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上で,

$$\arg \min_{f \in \mathcal{H}} \mathcal{L}_D(f)$$

を求める問題を単回帰分析という．単回帰分析の最適仮説 $f^* \in \mathcal{H}$ は

$$f^*(w) = \frac{\sum_{i=1}^N x_i y_i}{\sum_{i=1}^N x_i^2} w$$

となる．また, $\mathcal{H} \simeq \mathbb{R}$ である．

Example 3.2. (重回帰分析)

機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を以下のように定義する．

$$\mathcal{X} = \mathbb{R}^d (N \geq d), \mathcal{Y} = \mathbb{R},$$

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f(x) = W^\top x + b, W \in \mathbb{R}^d, b \in \mathbb{R}\},$$

$$\mathcal{L}_D(f) = \sum_{i=1}^N |f(x_i) - y_i|^2.$$

この機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上で

$$\arg \min_{f \in \mathcal{H}} \mathcal{L}_D(f)$$

を求める問題を重回帰分析という． $X \in \mathbb{R}^{N \times (d+1)}$ を

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1d} \\ 1 & x_{21} & x_{22} & \cdots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \cdots & x_{Nd} \end{pmatrix}$$

と定義し, $\mathbf{y} = (y_1, y_2, \dots, y_N)^\top \in \mathbb{R}^N$ とする. この時, $\text{rank}(X) = d + 1$ ならば, 最適仮説 $f^* \in \mathcal{H}$ は $[b \ W] = (X^\top X)^{-1} X^\top \mathbf{y}$ とすれば

$$f^*(x) = [b \ W]^\top x$$

である.

3.2 線形回帰

前小節で定義した重回帰分析は線形的なデータにしか適応できなかった. それを基底関数と呼ばれる関数列を定義する事で, 非線形データにも対応することができる.

Definition 3.3. (基底関数)

\mathcal{X} を入力空間とし, 任意に $d \in \mathbb{N}$ をとる. \mathcal{X} 上の連続関数列 $\{\phi_n\}_{n=1}^d \subset C(\mathcal{X}, \mathbb{R})$ が一次独立であるとき, $\Phi = (\phi_1, \phi_2, \dots, \phi_d)^\top : \mathcal{X} \rightarrow \mathbb{R}^d$ を基底関数 (basis function) と呼ぶ.

Example 3.4. (多項式回帰)

機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を以下のように定義する.

$\mathcal{X} = \mathbb{R}, \mathcal{Y} = \mathbb{R},$

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f(x) = W^\top \Phi(x), W \in \mathbb{R}^{d+1}\},$$

$$\mathcal{L}_D(f) = \sum_{i=1}^N |f(x_i) - y_i|^2.$$

ここで, $\phi_n(x) = x^n, n \in \{0, 1, \dots, d\}$ とする.

この機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上で

$$\arg \min_{f \in \mathcal{H}} \mathcal{L}_D(f)$$

を求める問題を多項式回帰という. $X = [\phi_1(x_1), \phi_2(x_2), \dots, \phi_N(x_N)]^\top \in \mathbb{R}^{N \times d}$ とし, $\mathbf{y} = (y_1, y_2, \dots, y_N)^\top \in \mathbb{R}^N$ とする. この時, 最適仮説 f^* は重回帰分析と同様に $W = (X^\top X)^{-1} X^\top \mathbf{y}$ とすれば

$$f^*(x) = W^\top x$$

である.

Remark 3.5. (ハイパーパラメータと学習パラメータ)

多項式回帰の多項式の次数 $d \in \mathbb{N}$ や基底関数 $\{\phi_n\}_{n=1}^d$ のようにコンピュータに学習させるのではなく, 機械学習モデルの設計者が設定するパラメータのことをハイパーパラメータと呼ぶ. 一方, $W \in \mathbb{R}^{d+1}$ のように, データからコンピュータが自動で学習するパラメータのことを学習パラメータと呼ぶ.

3.3 過学習と正則化

前節の多項式回帰で学習を行うと以下のグラフのように、訓練データに過剰適合してしまい、未知のデータについての汎化性能が小さくなってしまいうことを過学習 (overfitting) という。

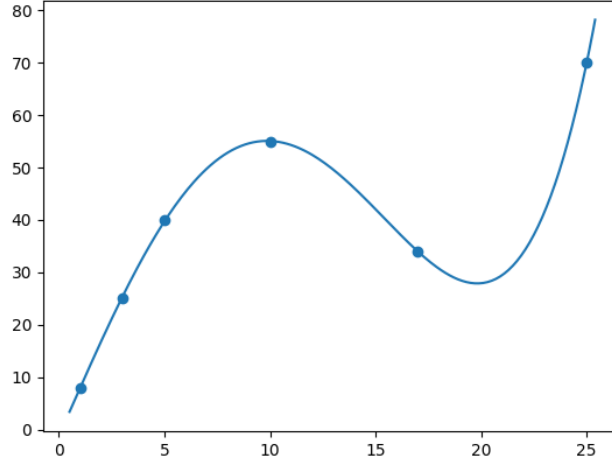


図1 $d = 4$ の時の多項式回帰の過学習

過学習を防ぐための手法が正則化である。

Definition 3.6. (正則化)

$\mathcal{L}_D : \mathcal{H} \rightarrow \mathbb{R}$ を損失関数とする. $\mathcal{L}_D^R := \mathcal{L}_D + L^R$ とする. この時, \mathcal{L}_D^R を \mathcal{L}_D の正則化と呼び, L^R を正則化項と呼ぶ.

Example 3.7. (Lasso 正則化多項式回帰)

$\lambda \in \mathbb{R}^+$, $d \in \mathbb{N}$ を任意にとる. 機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を以下のように定義する.

$\mathcal{X} = \mathbb{R}, \mathcal{Y} = \mathbb{R}$,

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f(x) = W^\top \Phi(x), W \in \mathbb{R}^{d+1}\},$$

$$\mathcal{L}_D(f) = \sum_{i=1}^N |f(x_i) - y_i|^2 + \lambda W^\top W. \quad (W \text{ は } f \text{ のパラメータ})$$

ここで, $\phi_n(x) = x^n$, $n \in \{0, 1, \dots, d\}$ とする.

この機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上で

$$\arg \min_{f \in \mathcal{H}} \mathcal{L}_D(f)$$

を求める問題を Lasso 正則化多項式回帰という. 通常多項式回帰の時と同様に, $X = [\phi_1(x_1), \phi_2(x_2), \dots, \phi_N(x_N)]^\top \in \mathbb{R}^{N \times d}$ とし, $\mathbf{y} = (y_1, y_2, \dots, y_N)^\top \in \mathbb{R}^N$ とする. この時,

最適仮説 f^* は $W = (X^\top X + \lambda I)^{-1} X^\top \mathbf{y}$ とすれば

$$f^*(x) = W^\top x$$

である.

正則化の結果が以下のグラフである.

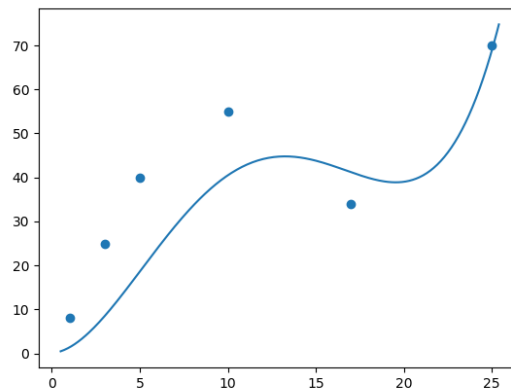


図2 $\lambda = 700$, $d = 4$ の時の Lasso 正則化多項式回帰

このように正則化をすることで, パラメータが大きくなりすぎることを防ぐことで過学習を防ぐことができる.

Remark 3.8. (未学習について)

正則化を行う際にパラメータ $\lambda \in \mathbb{R}^{+*3} := \{x \in \mathbb{R} \mid x \geq 0\}$ は自分で決める必要がある. その際に λ の値を大きくしすぎると未学習 (underfitting) という問題が発生する. 未学習とは訓練データに対してモデルが十分に適合できていないことをいう.

*3 正則化パラメータという.

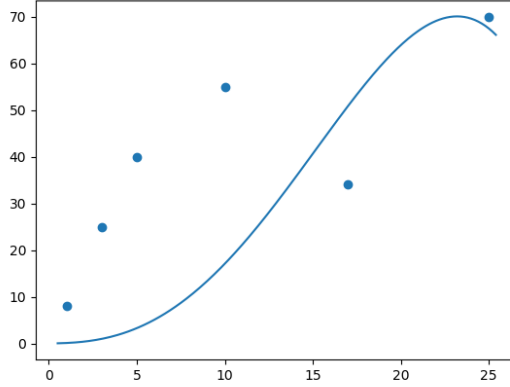


図3 $\lambda = 30000$, $d = 4$ の時の Lasso 正則化多項式回帰の未学習

このように, 正則化パラメータは適切に選ぶことが大切である.

3.4 ロジスティック回帰と勾配降下法

最後の具体例として, ロジスティック回帰を紹介する. ロジスティック回帰は分類問題を解くための手法の 1 つである. まず, 実数値ベクトルを確率に変換するソフトマックス関数および, one-hot ベクトルというものを導入する.

Definition 3.9. (ソフトマックス関数)

$\psi : \mathbb{R}^d \rightarrow \mathbb{R}^d$ を以下で定義する.

$$\psi(x)_i = \frac{\exp(x_i)}{\sum_{k=1}^d \exp(x_k)}$$

この時, ψ をソフトマックス関数と呼ぶ.

Definition 3.10. (one-hot ベクトル)

$y \in \mathbb{R}^m$ が第 $c \in \{1, 2, \dots, m\}$ 成分が 1 であり, 残りの成分が 0 である時, y を one-hot ベクトルという.

Example 3.11. (m 値分類ロジスティック回帰)

機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ を以下のように定義する.

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}^m$,

$$\mathcal{H} = \{f : \mathcal{X} \rightarrow \mathcal{Y} \mid f(x) = \psi(Wx + b), W \in \mathbb{R}^{m \times d}, b \in \mathbb{R}^m\},$$

$$\mathcal{L}_D(f) = - \sum_{n=1}^N \sum_{k=1}^m y_{nk} \log f(x_n)_k,$$

ここで $\psi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ はソフトマックス関数であり, y_n は one-hot ベクトルである.
この機械学習空間 $(\mathcal{X}, \mathcal{Y}, D, \mathcal{H}, \mathcal{L}_D)$ 上で

$$\arg \min_{f \in \mathcal{H}} \mathcal{L}_D(f)$$

を求める問題を m 値分類ロジスティック回帰という.

先ほどまでは, 損失関数の最小値を解析的に求めているが, 今回の損失関数の最小値を解析的に求めることはせず, 勾配降下法という連続最適化アルゴリズムを用いて求めることにする. 勾配降下法とは C^∞ 関数 F の最小値を求めるためのアルゴリズムである.

Algorithm 1 Gradient Decent

Require: F : smooth function on \mathbb{R}^d

Require: $0 < \alpha < 1$: learning rate

Require: θ : Initial parameter vector

while θ not converged **do**

$\theta \leftarrow \theta - \alpha \nabla F(\theta)$

end while

return θ

勾配降下法の最適解への収束条件を述べておく.

Definition 3.12. (β -Lipschitz 連続)

$(X, d_X), (Y, d_Y)$ を距離空間とする. $f : X \rightarrow Y$ が β -Lipschitz 連続であるとは,

$$\forall x, y \in X, d_Y(f(x), f(y)) \leq \beta d_X(x, y)$$

が成立することである.

Theorem 3.13. (勾配降下法の最適化収束条件)

\mathcal{H} をヒルベルト空間とし, $f : \mathcal{H} \rightarrow \mathbb{R}$ を C^∞ 関数とする. この時, ∇f が β -Lipschitz 連続ならば, 学習率 $\alpha \in \left(0, \frac{2}{\beta}\right)$ で最適解に弱収束する.

Example 3.14. (2 次関数における勾配降下法)

$y = x^2$ において勾配降下法を適応した結果を以下の図と表に示す. 本来の最適値である $x = 0$ に近づいているのが分かる.

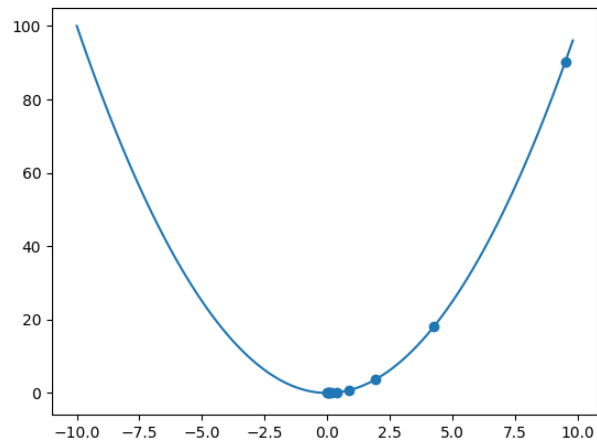


図4 $\alpha = 0.001$ とした時の勾配降下法の結果.

イタレーション数	0	400	800	1200	1600	2000	2400	2800	3200	3600
θ の値	9.50	4.27	1.91	0.86	0.39	0.17	0.08	0.03	0.02	0.01

表1 値の詳細

3.5 MNIST での実験

MNIST とは下記のような 0 から 9 までの自然数が書かれた画像のデータセットの名前である. 今回は, m 値分類のロジスティック回帰を用いて、画像に書かれている文字を予測する問題を考える. ロジスティック回帰の損失関数 \mathcal{L}_D を勾配降下法を用いて最小化するため, まずは損失関数の勾配を計算する. ここで, $\mathcal{H} \simeq \mathbb{R}^{m \times d+1}$ だから, \mathcal{L}_D を $\mathbb{R}^{m \times d+1}$ 上の関数と考える.

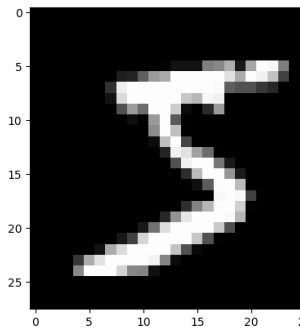


図5 数字の 5

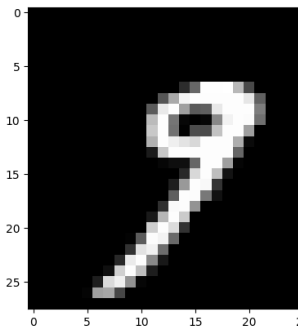


図6 数字の 9

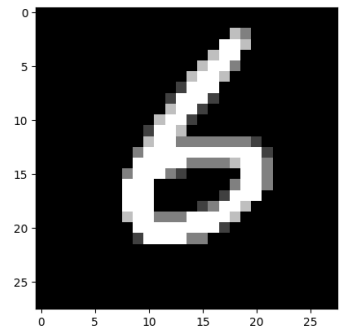


図7 数字の 6

Lemma 3.15. 任意の $i \in \{1, 2, \dots, m\}, j \in \{1, 2, \dots, d\}$ に対して以下が成立する.

$$\begin{aligned}\frac{\partial \mathcal{L}_D}{\partial w_{ij}} &= - \sum_{n=1}^N (y_{ni} - f(x_n)_i) x_{nj} \\ \frac{\partial \mathcal{L}_D}{\partial b_i} &= - \sum_{n=1}^N (y_{ni} - f(x_n)_i)\end{aligned}$$

この補題より,

$$\begin{aligned}\frac{\partial \mathcal{L}_D}{\partial W} &= - \begin{pmatrix} \sum_{n=1}^N (y_{n1} - f(x_n)_1) x_{n1} & \cdots & \sum_{n=1}^N (y_{n1} - f(x_n)_1) x_{nd} \\ \vdots & \ddots & \vdots \\ \sum_{n=1}^N (y_{nm} - f(x_n)_m) x_{n1} & \cdots & \sum_{n=1}^N (y_{nm} - f(x_n)_m) x_{nd} \end{pmatrix} \\ &= - \begin{pmatrix} y_{11} - f(x_1)_1 & \cdots & y_{N1} - f(x_N)_1 \\ \vdots & \ddots & \vdots \\ y_{1m} - f(x_1)_m & \cdots & y_{Nm} - f(x_N)_m \end{pmatrix} \begin{pmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{N1} & \cdots & x_{Nd} \end{pmatrix} \\ &= (Y - F)X^\top,\end{aligned}$$

ここで, $X = [x_1, x_2, \dots, x_N]$, $Y = [y_1, y_2, \dots, y_N]$, $F = [f(x_1), f(x_2), \dots, f(x_N)]$ である. この補題より, また, この計算より $\mathbf{1} = [1, 1, \dots, 1]^\top \in \mathbb{R}^{N \times 1}$ とすれば,

$$\frac{\partial \mathcal{L}_D}{\partial b} = (Y - F)\mathbf{1}$$

もいえる. これらを用いて学習した結果を以下に示す.

4 Deep Learning

4.1 Neural Network

4.2 The Universal Theorem of Neural Network

5 Generative Adversarial Networks

5.1 GAN の定式化

$(\Omega, \mathcal{F}, \mathbb{P})$ を完備な確率空間とする. \mathcal{X} を d 次元線型空間とする.

$$\begin{aligned}\mathcal{H}_1 &= \{G: \mathcal{Z} \rightarrow \mathcal{X} \mid G \text{ はニューラルネット} \}, \\ \mathcal{H}_2 &= \{D: \mathcal{X} \rightarrow [0, 1] \mid D \text{ はニューラルネット} \}\end{aligned}$$

ここで, \mathcal{Z} は潜在空間と呼ばれる \mathbb{R}^d の線型部分空間である. また, 確率変数 $Z: \Omega \rightarrow \mathcal{Z}$ に対し, $g(Z)$ が従う確率分布を \mathbb{P}

6 Applications of GANs

参考文献

- [1] Ian J. 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*, 2014.
- [2] Martin Arjovsky and Leon Bottou, Towards Principled Methods for Training Generative Adversarial Networks, *International Conference on Learning Representations*, 2017.
- [3] Takeru Miyato, Toshiki Kataoka, Masanori Koyama and Yuichi Yoshida, Spectral Normalization for Generative Adversarial Networks, *International Conference on Learning Representations*, 2018.
- [4] 矢田部 浩平, 信号処理と最適化, http://www.sp.ipc.i.u-tokyo.ac.jp/~saruwatari/SP-Grad2018_09ProfYatabe.pdf, 早稲田大学 表現工学科, 2018.
- [5] <http://darden.hatenablog.com/entry/2018/01/25/222201>