

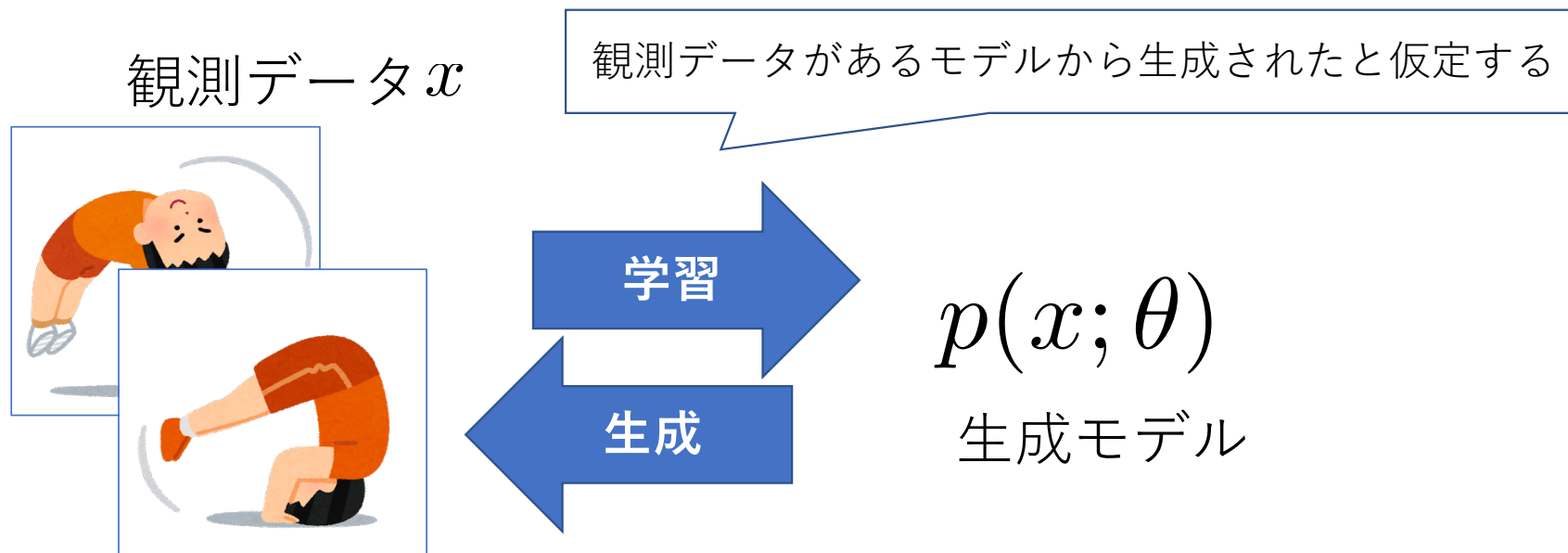
GANとは

5/18/2018

sakai

- 生成モデル
- GAN
- GANの欠点
- 色々なGAN

生成モデル

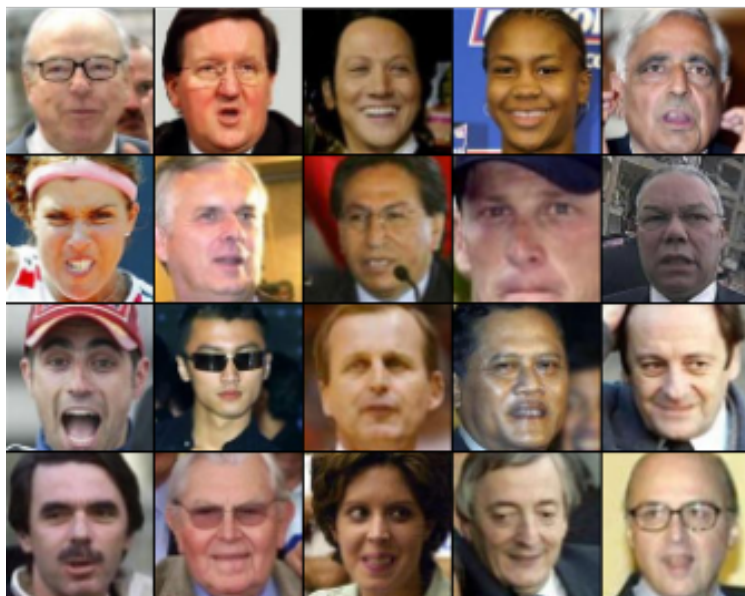


最近の流行りはこの生成モデルを
NNによって表現すること。

$$z \rightarrow \boxed{\text{生成器}} \rightarrow x$$
$$x = f(z)$$

- サンプリングが可能.
 - 確率モデルから未知のデータの生成.
- データに対する確率密度が得られる.
- 半教師あり学習に応用ができる.

Input



VAE reconstruction



- 尤度の評価のため単純な分布(ガウス, ベルヌーイ)を仮定していることが原因(多分).

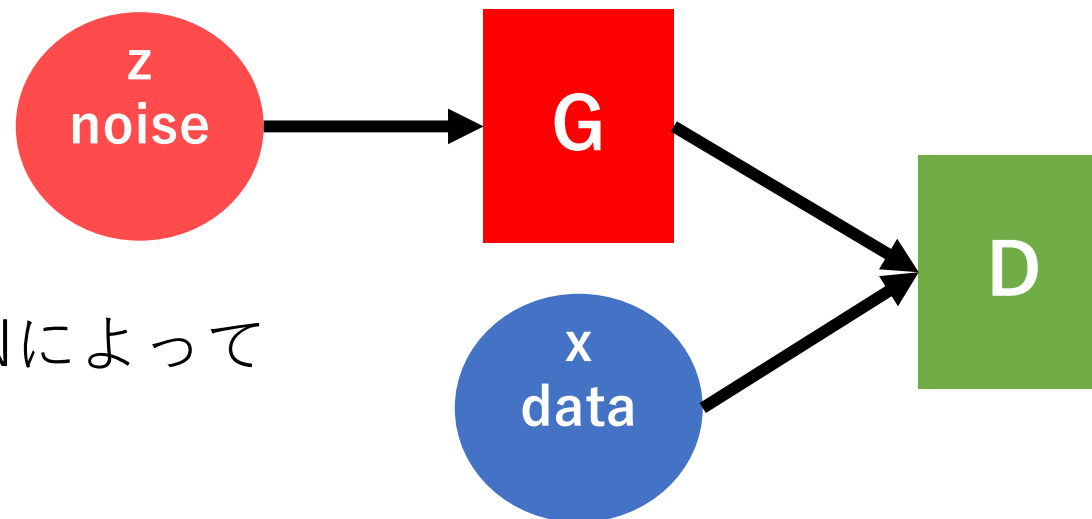
GAN

パラメトリックな分布を陽に定義する.

- ▶ パラメータをNNで得る.
- ▶ 尤度最大化でデータ分布にモデル分布を近づける.

分布を仮定せず
分布の形そのものをNNで表現しよう.

- 訓練データ $x \sim p_{data}(x)$ を表現する生成器Gを得ることが目的.
- つまり生成器Gの出力の分布 $p_g(x)$ を $p_{data}(x)$ に近づけることが目的である.



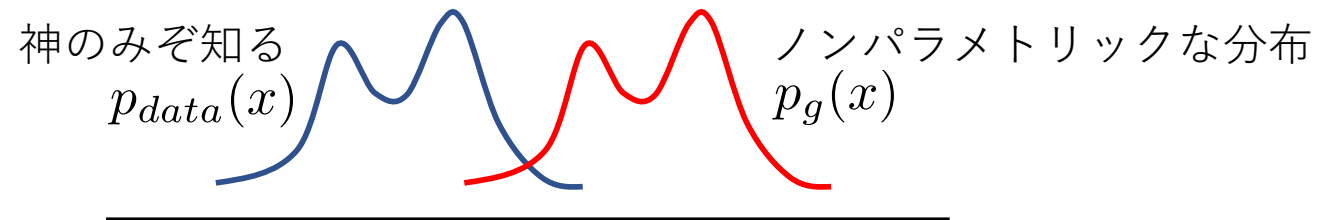
敵対する二つのNNによって
訓練を行うモデル.

ただし $p_g(x)$ は未定義のままです。

- ▶ 未定義の分布に対して尤度は測れない

目標 データの分布 $p_{data}(x)$ に近いモデル分布 $p_g(x)$ を求める.

モデル分布: 生成器Gの出力が従う分布



- ▶ 分布の比較ができればいいので密度比を考える.

$$r(x) = \frac{p_{data}(x)}{p_g(x)}$$

このままではまだカーネル密度推定などで求める必要がある.

- データの分布に従うデータと生成器から生成したデータを同じ割合で用意し、ラベル付けされたデータ集合を考える。
- データの分布に従うデータを $y=1$, 生成器から生成されたデータを $y=0$ とした。

$$\mathcal{D} = \{(x_1, 0), \dots, (x_N, 0), (x_{N+1}, 1), \dots, (x_{2N}, 1)\}$$

$$\text{ただし, } p(y=1) = p(y=0) = \frac{1}{2}$$

このとき, それぞれの分布は

$$p_{data}(x) = p(x|y=1) \quad : \text{データの従う分布}$$

$$p_g(x) = p(x|y=0) \quad : \text{生成器の出力値が従う分布}$$

密度比からクラス分類問題へ

12/53

密度比は以下のようなになる。

$$\frac{p_{data}(x)}{p_g(x)} = \frac{p(x|y=1)}{p(x|y=0)} = \frac{\frac{p(y=1|x)p(x)}{p(y=1)}}{\frac{p(y=0|x)p(x)}{p(y=0)}} = \frac{p(y=1|x)}{p(y=0|x)}.$$

$p(y=1|x)$ を学習で推定したい。

NNで上記を推定する分布をパラメータ化する。

$$x \rightarrow \boxed{D(x; \theta_d)} \rightarrow y$$

識別器
discriminator
 $D(x; \theta_d) = p(y=1|x)$

密度推定を二値分類問題に置き換えた。

$$p(y|x) = \text{Bern}(y|D(x; \theta_d))$$

➤ 識別器のパラメータ θ_d を決定する.

$$\begin{aligned} \text{尤度: } \quad p(\mathbf{y}|D(\mathbf{x}; \theta_d)) &= \prod_{i=1}^{2N} D(x_i; \theta_d)^{y_i} (1 - D(x_i; \theta_d))^{1-y_i} \\ \text{where } \quad \mathbf{y} &= (y_1, \dots, y_{2N}) \quad \mathbf{x} = (x_1, \dots, x_{2N}) \end{aligned}$$

➤ 次の負の交差エントロピーを最大化.

$$V(D) = \mathbb{E}_{p(x,y)} [y \log(D(x; \theta_d)) + (1 - y) \log(1 - D(x; \theta_d))]$$

➤ 式変形.

$$\begin{aligned} V(D) &= \mathbb{E}_{p(x|y)p(y)} [y \log(D(x; \theta_d)) + (1 - y) \log(1 - D(x; \theta_d))] \\ &= \mathbb{E}_{p(x|y=1)p(y=1)} [\log(D(x; \theta_d))] + \mathbb{E}_{p(x|y=0)p(y=0)} [\log(1 - D(x; \theta_d))] \\ &= \frac{1}{2} \mathbb{E}_{p_{data}(x)} [\log(D(x; \theta_d))] + \frac{1}{2} \mathbb{E}_{p_g(x)} [\log(1 - D(x; \theta_d))] \end{aligned}$$

➤ 定数倍は無視して,

$$V(D) = \mathbb{E}_{p_{data}(x)} [\log(D(x; \theta_d))] + \mathbb{E}_{p_g(x)} [\log(1 - D(x; \theta_d))]$$

目標 データの分布 $p_{data}(x)$ に近いモデル分布 $p_g(x)$ を求める.

識別関数が適切に推定できたなら

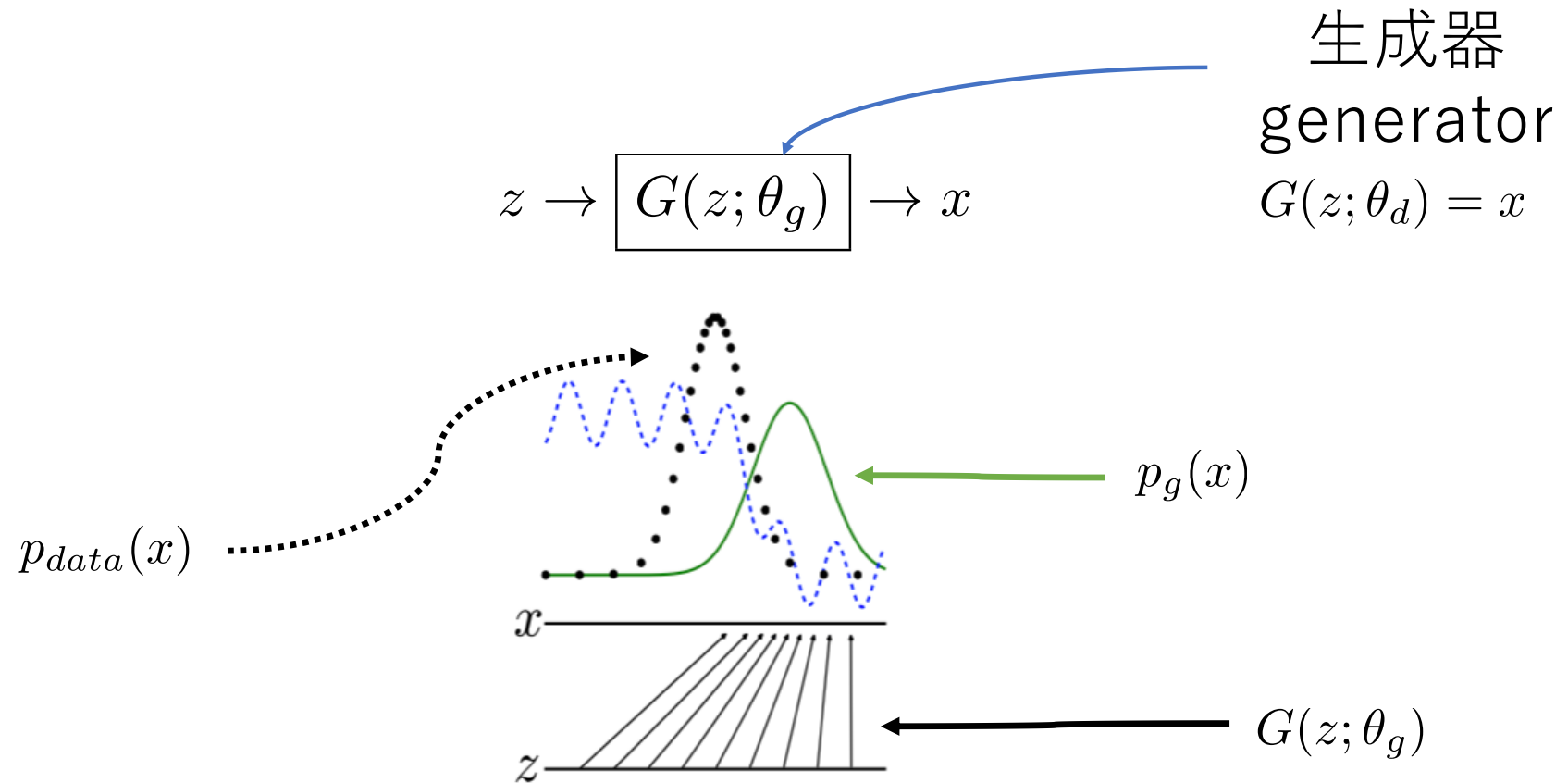
$$D^*(x; \theta_d^*) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \quad \text{となる(証明略).}$$

このとき, 目的関数は

$$V(D^*) = 2 \cdot JSD(p_{data} || p_g) - \log 4 \quad \text{となる(証明略).}$$

- JSDはJensen-Shannon divergence.
- Kullback-Leiber divergence と違い対称性がある.
- 識別器Dによって生成器Gのための評価指標(JSD)を用意した.

$p_{data}(x)$ を推定する分布をNNでパラメータ化する。



目標 データの分布 $p_{data}(x)$ に近いモデル分布 $p_g(x)$ を求める。

➤ $JSD(p_{data}||p_g)$ を最小にする生成器のパラメータを求める。

$$\begin{aligned} JSD(p_{data}||p_g) &= \frac{1}{2} \left\{ KL \left(p_{data} \parallel \frac{p_{data} + p_g}{2} \right) + KL \left(p_g \parallel \frac{p_{data} + p_g}{2} \right) \right\} \\ &= \mathbb{E}_{x \sim p_{data}} \left[\log \frac{2p_{data}(x)}{p_{data}(x) + p_g(x)} \right] + \mathbb{E}_{x \sim p_g} \left[\log \frac{2p_g(x)}{p_{data}(x) + p_g(x)} \right] \\ &= \mathbb{E}_{x \sim p_{data}} [\log D^*(x)] + \mathbb{E}_{x \sim p_g} [\log (1 - D^*(x))] - \log 4 \\ &= \mathbb{E}_{x \sim p_{data}} [\log D^*(x)] + \mathbb{E}_{z \sim p_z} [\log (1 - D^*(G(z)))] - \log 4 \end{aligned}$$

$$V(D, G) = \mathbb{E}_{p_{data}(x)} [\log(D(x; \theta_d))] + \mathbb{E}_{p_g(x)} [\log(1 - D(G(z; \theta_g); \theta_d))]$$

- 上記の式の最小化は以下の正答率の最小化とも解釈できる.

$$p(y = 1|\hat{x}) = D(\hat{x}; \theta_d) \quad \hat{x}: \text{真の分布から得られたデータ}$$

$$p(y = 0|G(z; \theta_g)) = 1 - D(G(z; \theta_g); \theta_d)$$

- GはDによってもたらされる近似的なJSD距離の最小化によって学習する.

実際には交互に目的関数を最適化する.

- Gを固定してDの学習.

$$\max_{\theta_d} \mathbb{E}_{p_{data}(x)} [\log(D(x; \theta_d))] + \mathbb{E}_{p_g(z)} [\log(1 - D(G(z; \theta_g); \theta_d))]$$

- Dを固定してGの学習.

$$\min_{\theta_g} \mathbb{E}_{p_g(z)} [\log(1 - D(G(z; \theta_g); \theta_d))]$$

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log \left(1 - D(G(z^{(i)})) \right) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

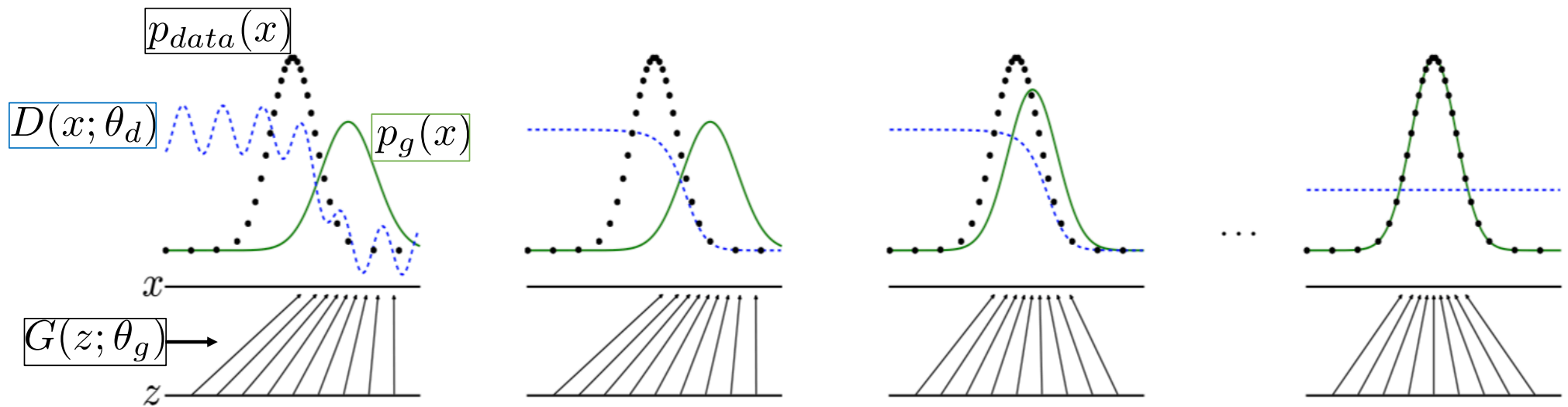
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(z^{(i)})) \right).$$

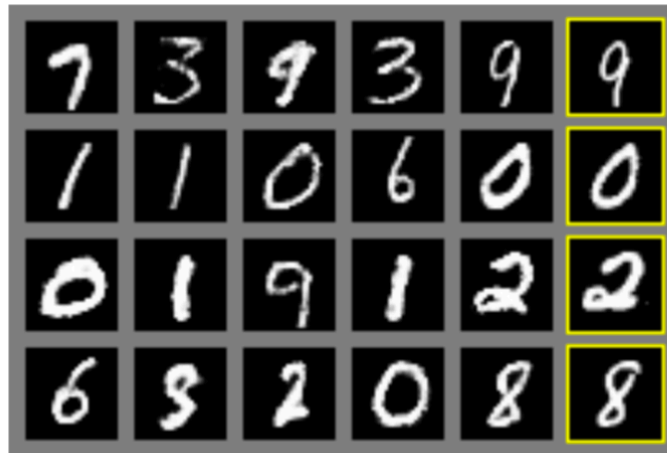
end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

GANの学習の様子

21/53

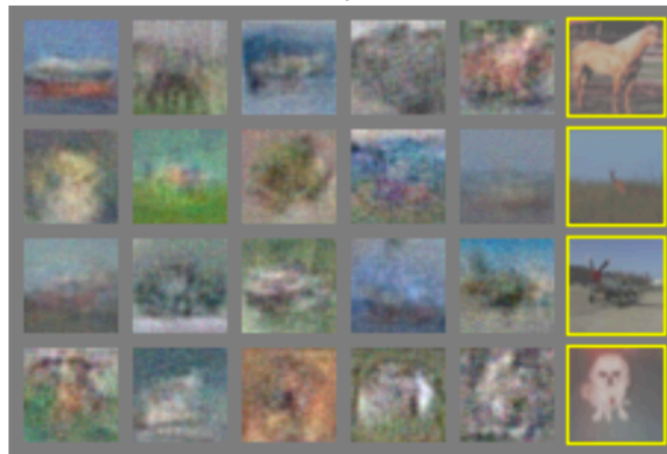




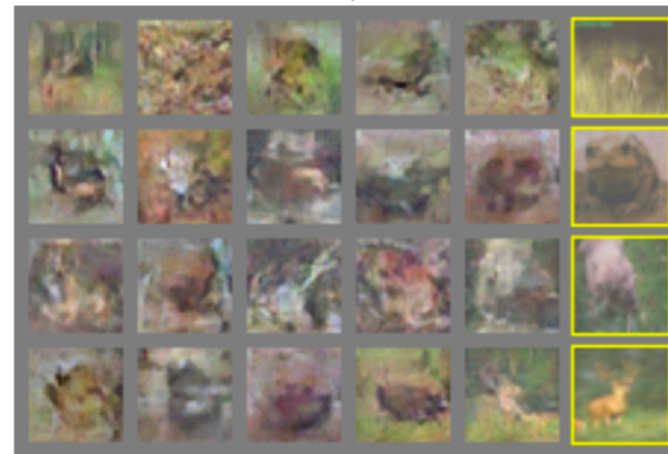
a)



b)



c)

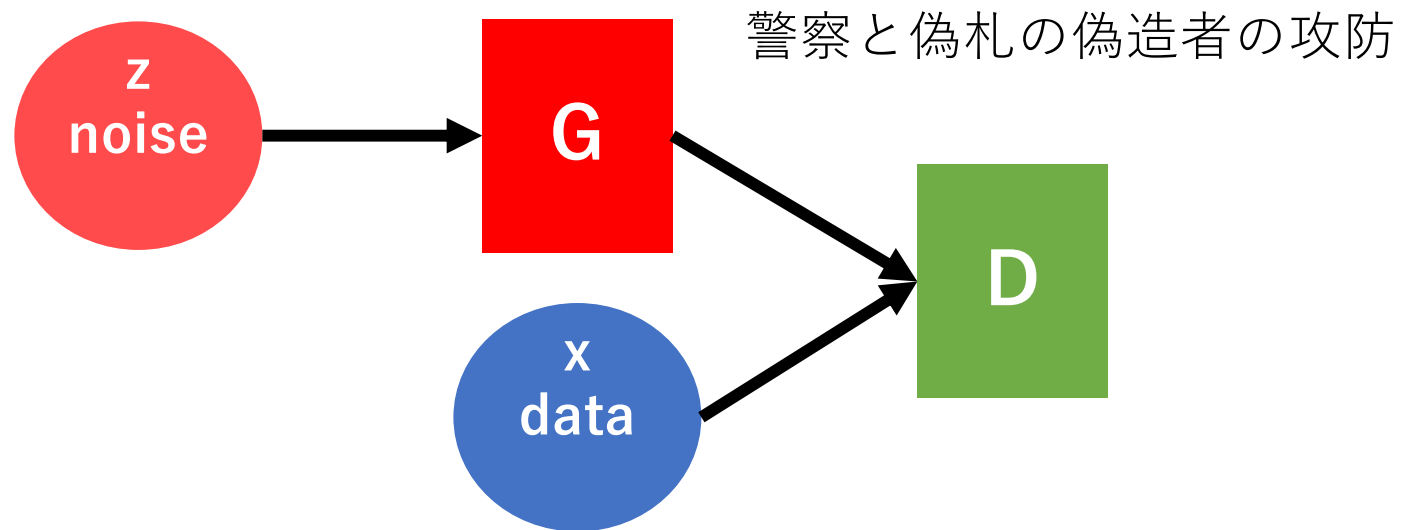


d)

直感的には以下のミニマックスゲームをする.

- GはDを騙すxを生成する.
- DはGに騙されないよう識別する.

$$\min_{\theta_g} \max_{\theta_d} V(D, G)$$



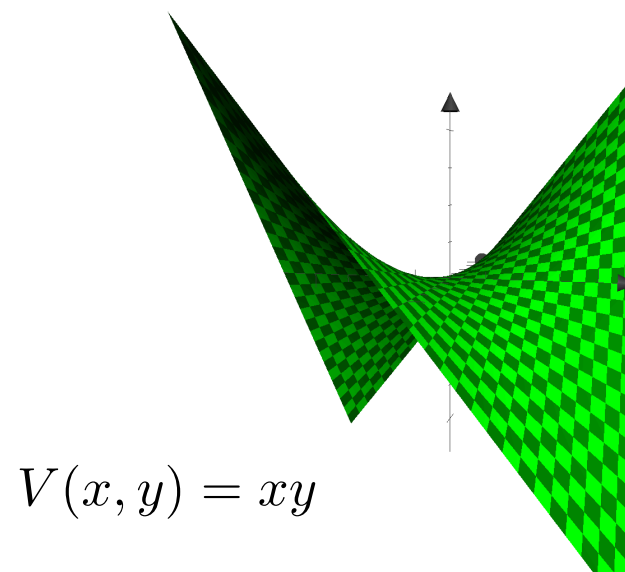
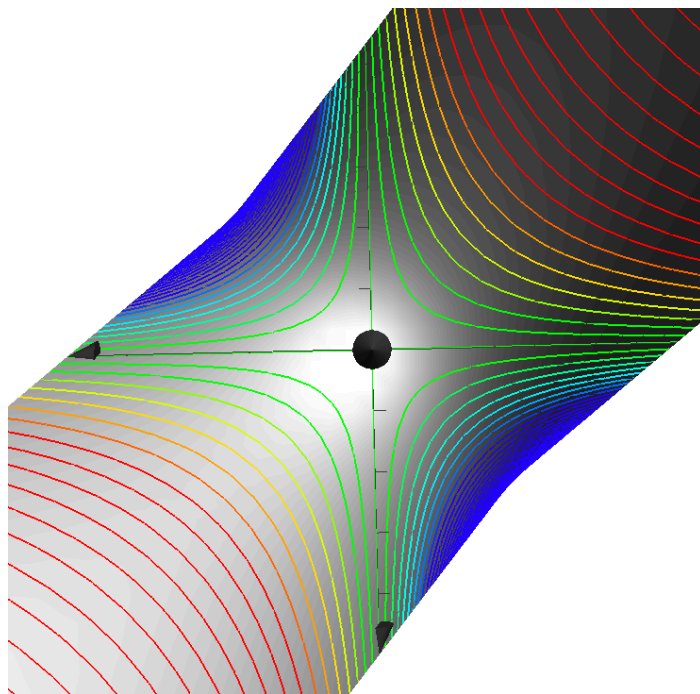
参考

Mohamed, Shakir, and Balaji Lakshminarayanan. "Learning in implicit generative models." *arXiv preprint arXiv:1610.03483*(2016).

GANの欠点

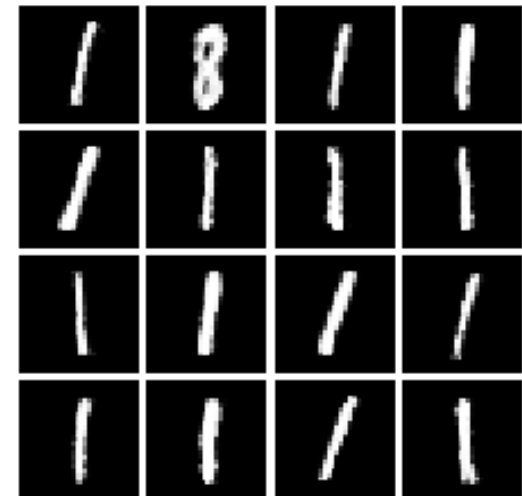
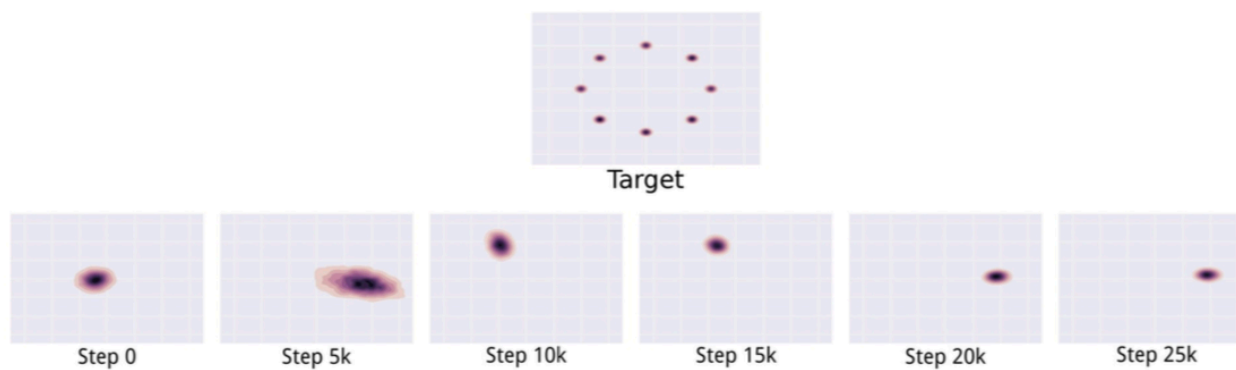
- 収束性
- Mode collapse
- 勾配消失

- ミニマックスゲームの平衡点.
 - 例えば $V(x, y) = xy$ 価値関数では $x = y = 0$ の鞍点が平衡点.
 - しかし x と y について交互に片方を固定し最適化すると振動する.

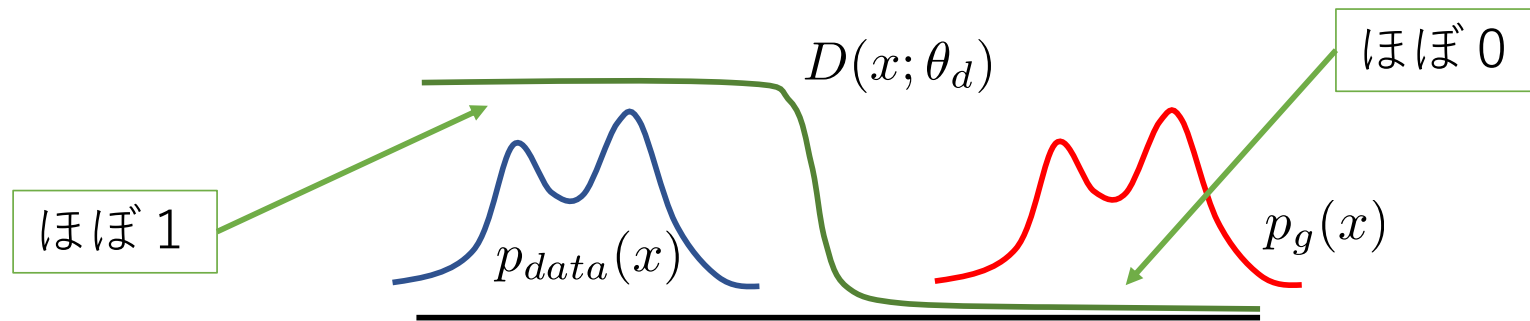


$$V(x, y) = xy$$

- 学習の不十分なDを固定してGを最適化すると…
 - Gの生成データが一つの峰にフィットする.



- 分布が交わらないとき(学習初期), 識別器が完全に識別できてしまう.



- 以下がサチってGについて更新ができなくなる.

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \left[\log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right) \right].$$

$$\min_G \mathbb{E}_{p_z} [\log (1 - D(G(z)))] \quad \longrightarrow \quad \min_G \mathbb{E}_{p_z} [-\log (D(G(z)))]$$

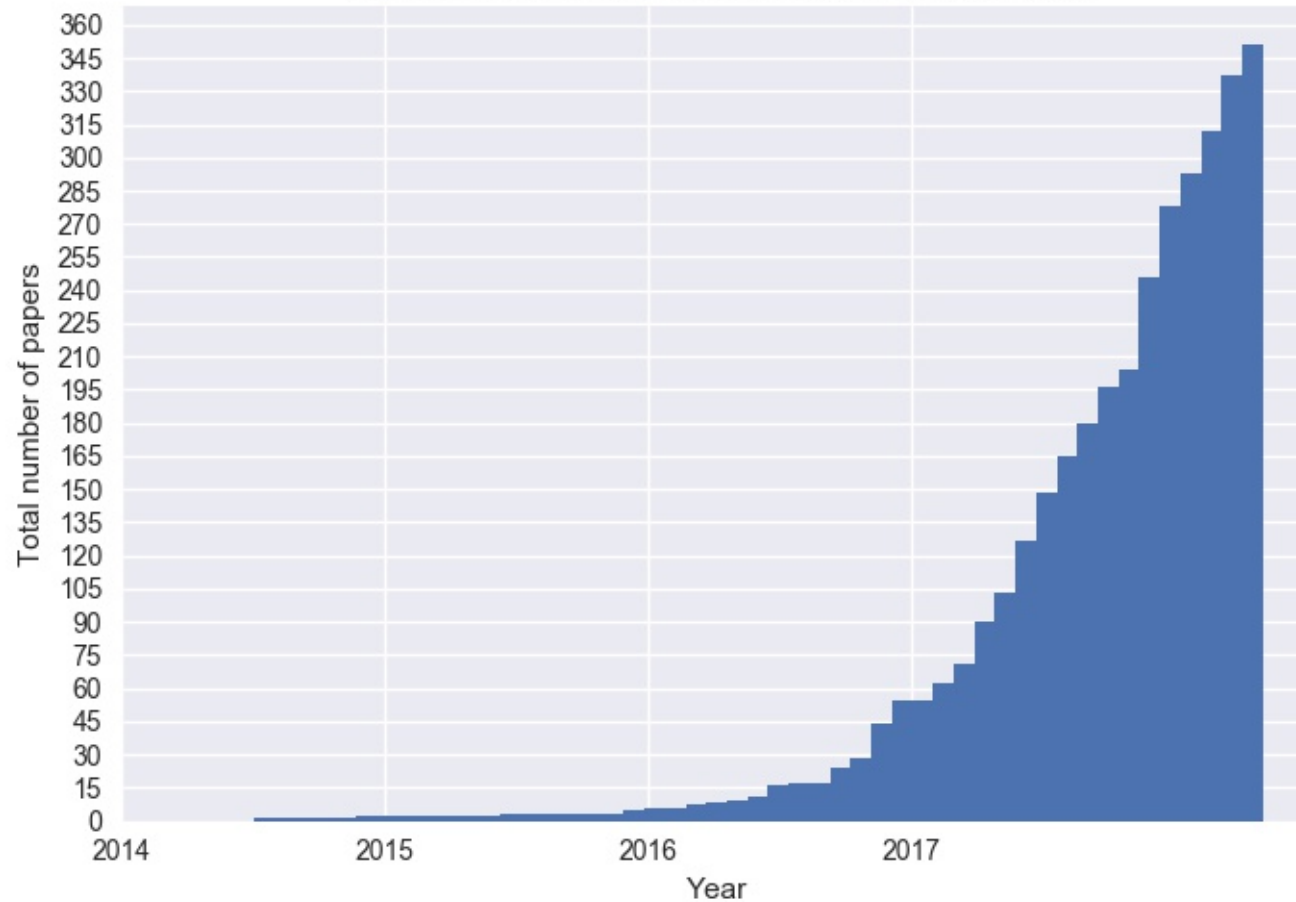
- Dがほぼ0を出力しても勾配がでる.
- もはやJSDとは関係がない.

- Dを完全に最適化してしまうと
 - $p_{data}(x)$ について1, $p_g(x)$ について0の定数を入力する.
 - 学習初期と同様Gについて勾配が消失する.

- GANの困りどころ
 - 学習が難しい.
 - Gが不完全なうちにDを最適化すると勾配が消失.
 - Dが不完全なうちにGを最適化するとmode collapse.

色々なGAN

Cumulative number of named GAN papers by month



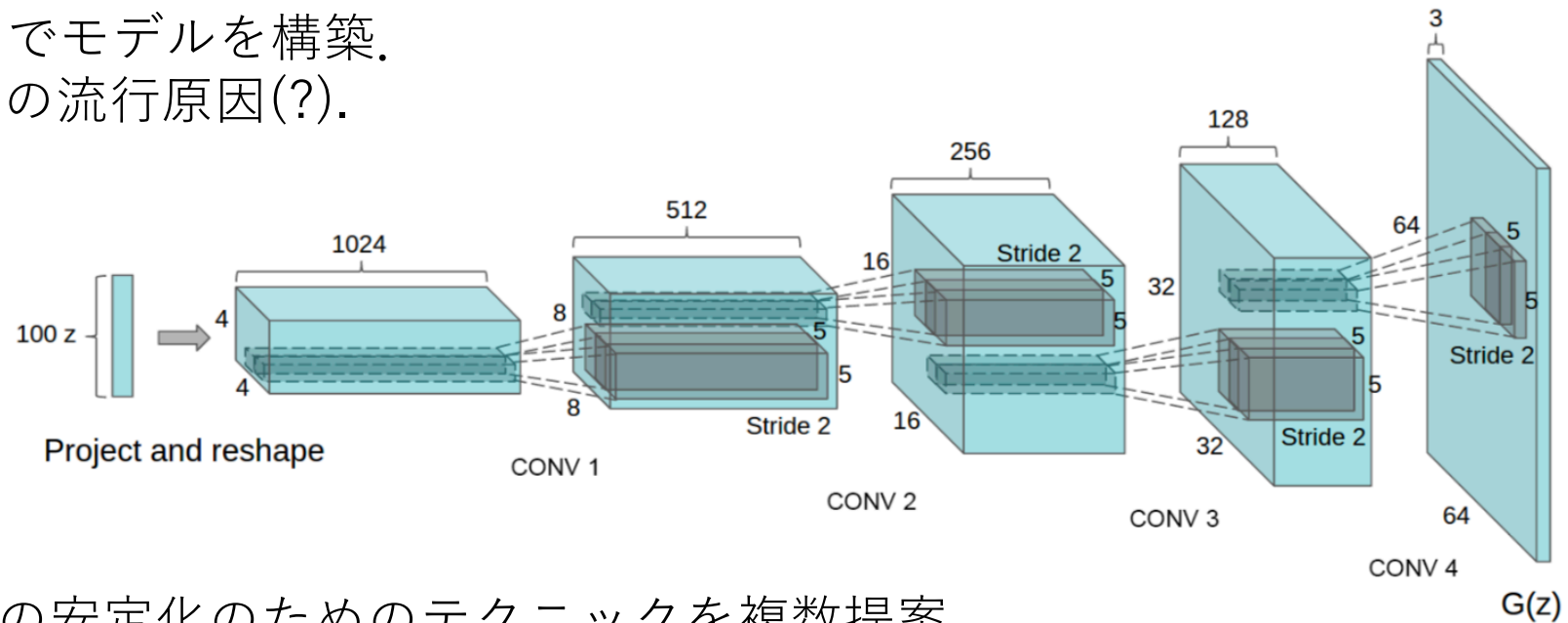
略称は被りがあるので
注意

[引用]GAN zoo
(2018/5/16)

(僕の中で)有名なGAN

- DCGAN(Deep Convolutional GAN)[2015]
- CGAN(Conditional GAN)[2014]
- Progressive GAN[2017]
- WGAN(Wasserstein GAN)[2017]

- CNNでモデルを構築.
- GANの流行原因(?).

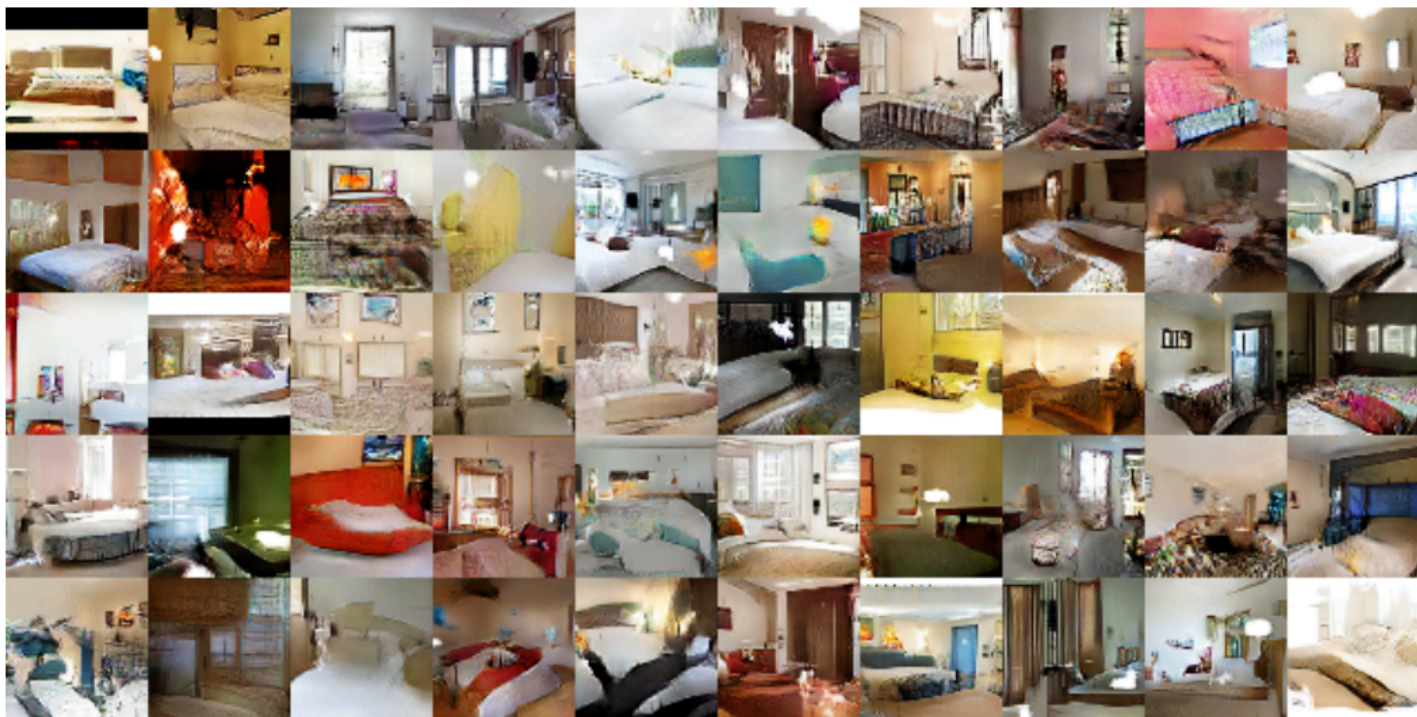


- 学習の安定化のためのテクニックを複数提案.
 - Batch Normalization.
 - Dにleaky ReLUを用いる.

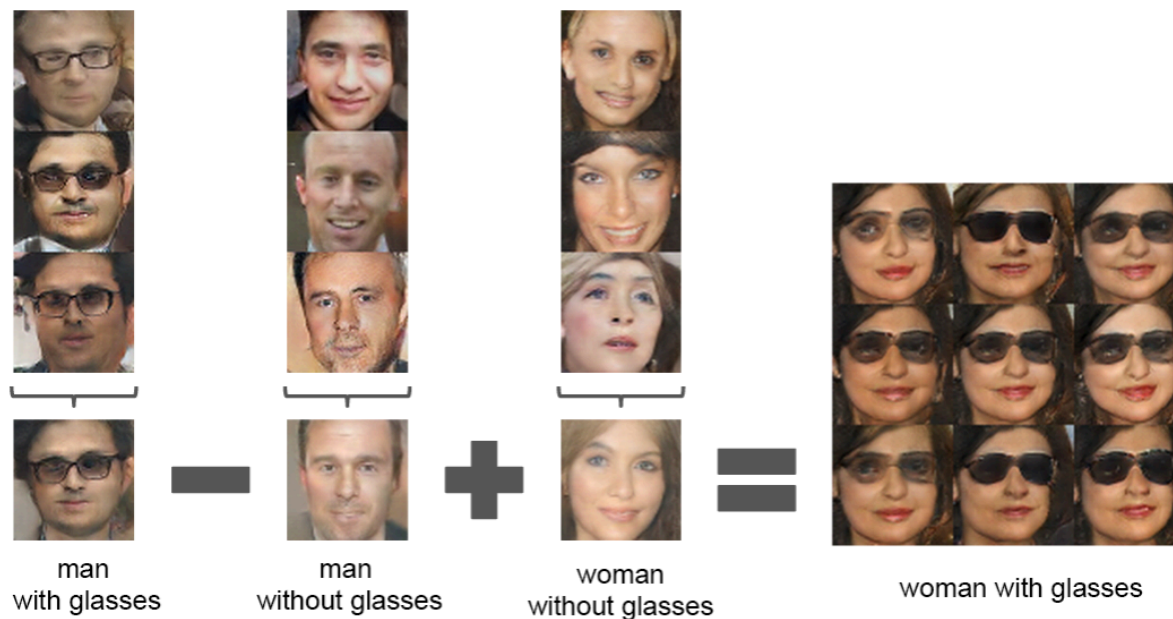
Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.

DCGANの生成物

35/53



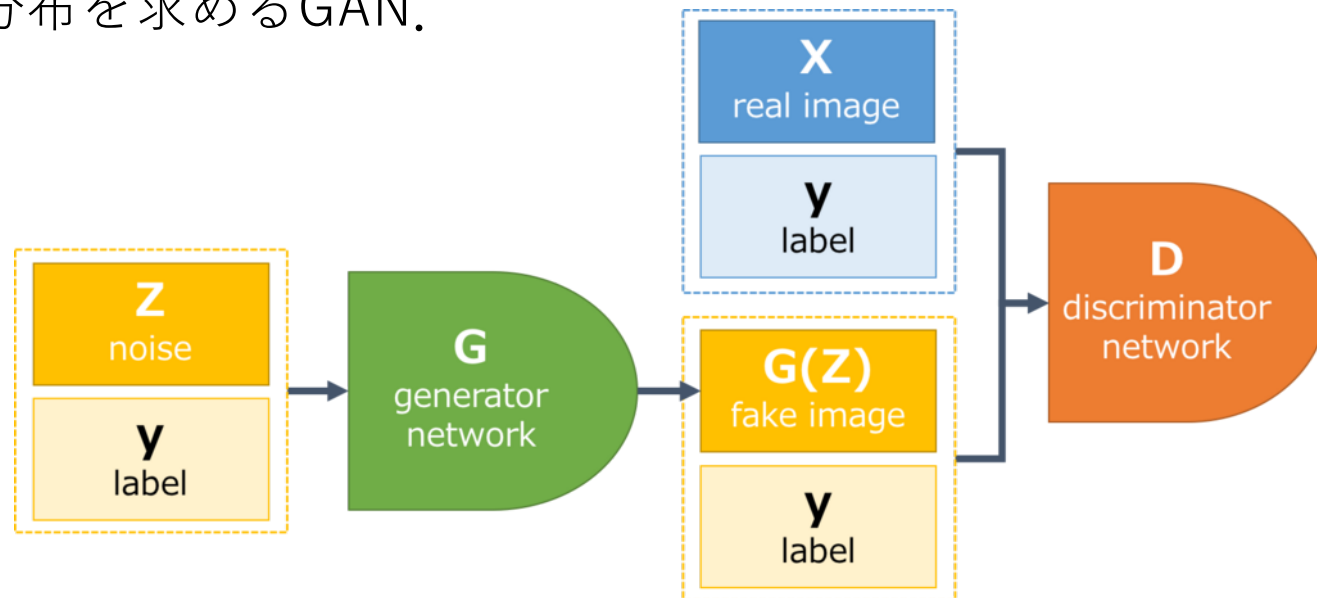
潜在空間での
演算



入力空間での
演算



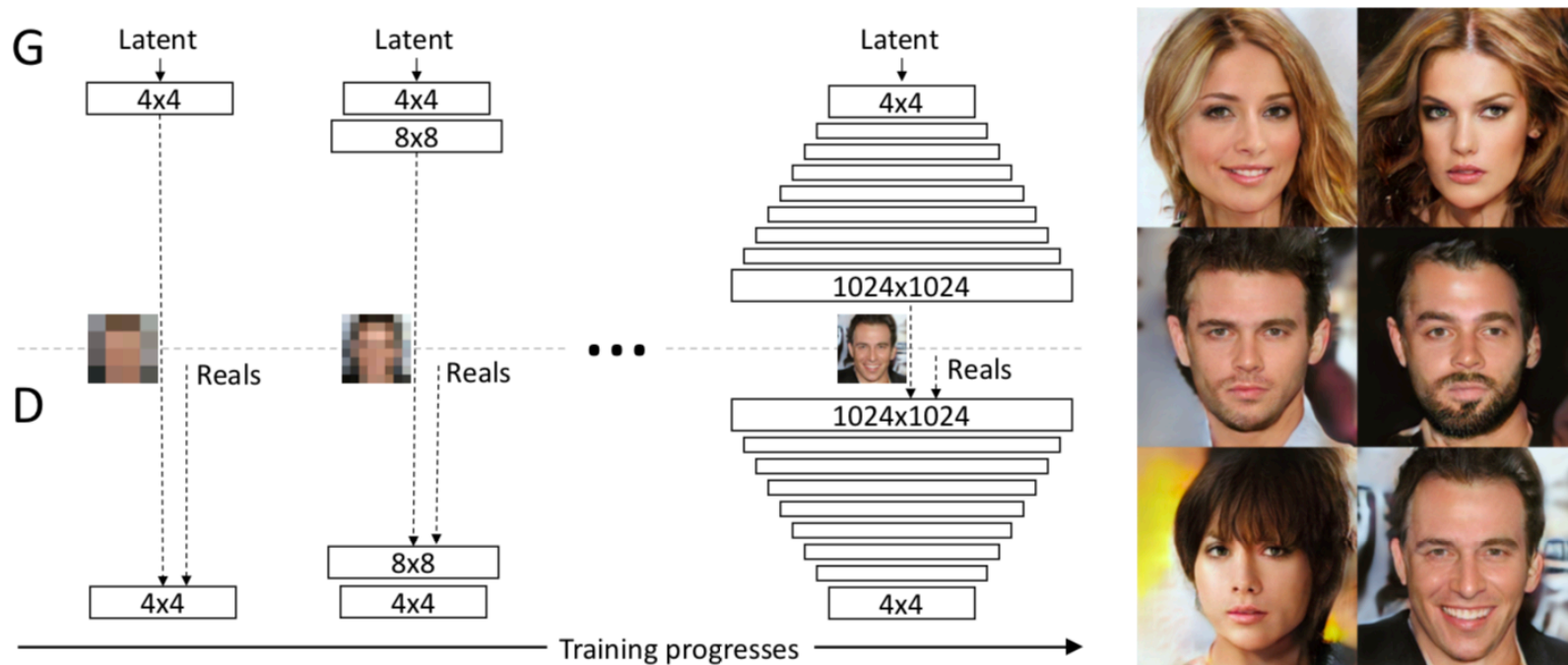
➤ 条件付き分布を求めるGAN.



$$\min_G \max_D V(D, G) = \mathbb{E}_{p_{data}(x)} [\log(D(x|y))] + \mathbb{E}_{p_g(x)} [\log(1 - D(G(z|y)))]$$

Mirza, M., & Osindero, S. (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.

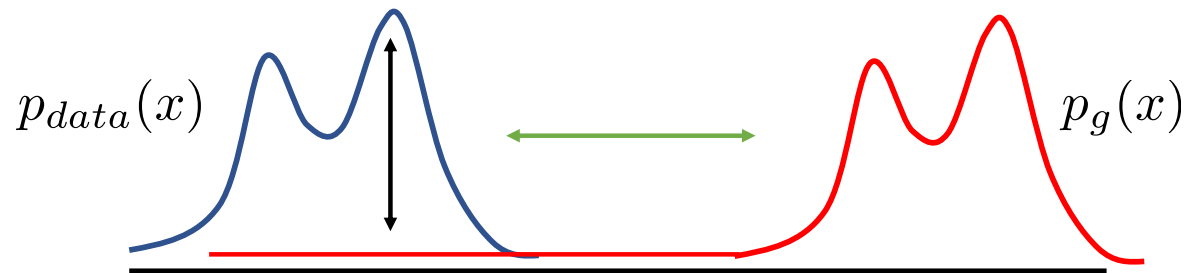
- 段階的に解像度を上げるように学習し，高解像度の画像生成に成功。



Karras, T., Aila, T., Laine, S., & Lehtinen, J. (2017). Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*.

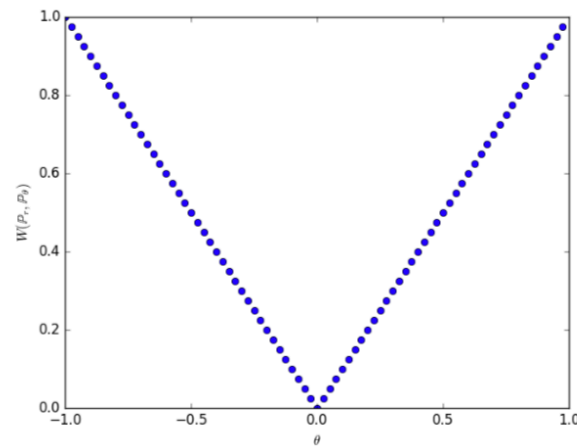
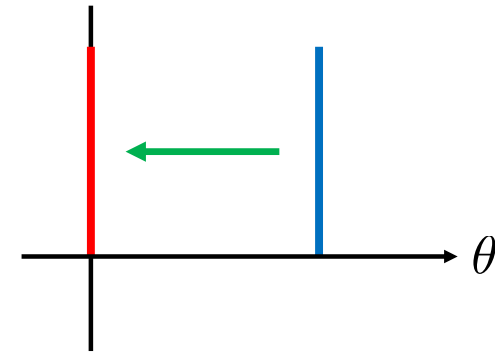
- 分布が交わらないとき，JSDは定数になる.

$$\begin{aligned}
 JSD(p_{data} || p_g) &= \frac{1}{2} \left\{ KL \left(p_{data} || \frac{p_{data} + p_g}{2} \right) + KL \left(p_g || \frac{p_{data} + p_g}{2} \right) \right\} \\
 &= \int p_{data}(x) \log \left\{ \frac{2p_{data}(x)}{p_{data}(x) + p_g(x)} \right\} dx + \int p_g(x) \log \left\{ \frac{2p_g(x)}{p_{data}(x) + p_g(x)} \right\} dx
 \end{aligned}$$

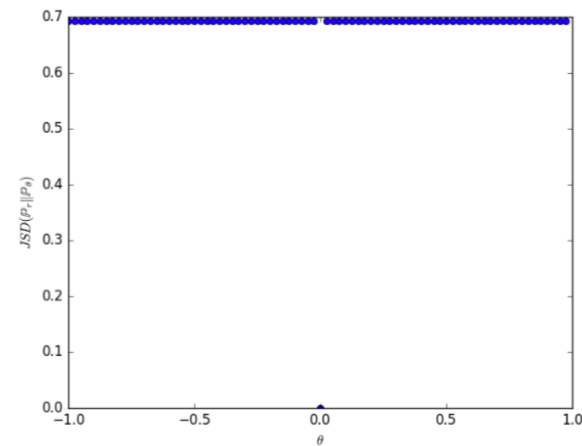


- 違うダイバージェンスを用いる提案が多くされている.
 - EMD(Earth Mover distance)

- 右の確率密度関数を動かした際の距離が以下.
- JSDでは広い範囲で定数, しかも不連続だが EMDは勾配が得られている.



EMD



JSD

- Wasserstein距離(の双対表現)による目的関数は以下.

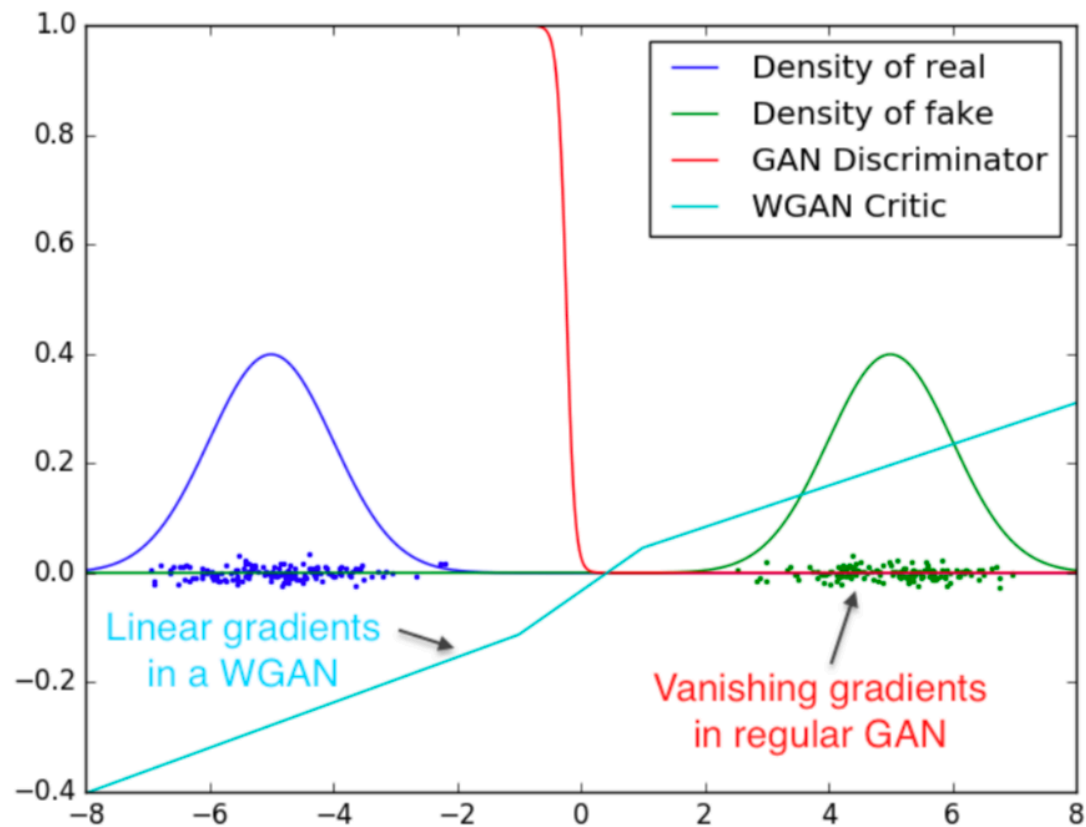
$$W(\mathbb{P}_r, \mathbb{P}_g) = \max_w \mathbb{E}_{x \sim \mathbb{P}_r} [D(x; w)] - \mathbb{E}_{z \sim p(z)} [D(G(z); w)].$$

ただし $\|D(x) - D(y)\| \leq \|x - y\|$ (1-リプシッツ連続)

- ほぼGANと同じ.
 - Dはリプシッツ性をもつ(傾きが急でない).
 - Dは非線型関数を出力に取らない(出力値は確率としての意味を持たない).
 - 付録.

Arjovsky, M., Chintala, S., & Bottou, L. (2017).
Wasserstein gan. *arXiv preprint arXiv:1701.07875*.

- Critic(GANでのDの出力に相当)は全域で勾配を与えておりGの学習を助ける.
 - Criticは確率ではない.



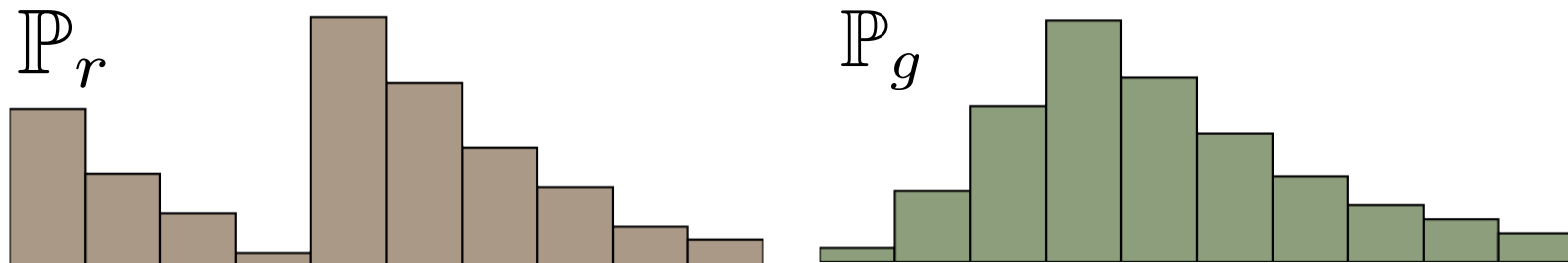
終わり

付録, 参考文献

- Earth-Mover distance or Wasserstein-1(以下EMD).

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|].$$

- 距離とは言うものの、それ自体が最適化問題.
 - 確率密度分布を土の山としてみたとき、片方の確率分布の山をもう一方の確率分布の山に土を輸送して変形させるための最小の労力.



- 山を目的の形へ変形させる輸送方法は無限に存在する.
- 輸送方法 $\gamma(x, y)$ を以下を満たすように定義する.

$$\sum_x \gamma(x, y) = \mathbb{P}_r(y) \quad \sum_y \gamma(x, y) = \mathbb{P}_g(x)$$

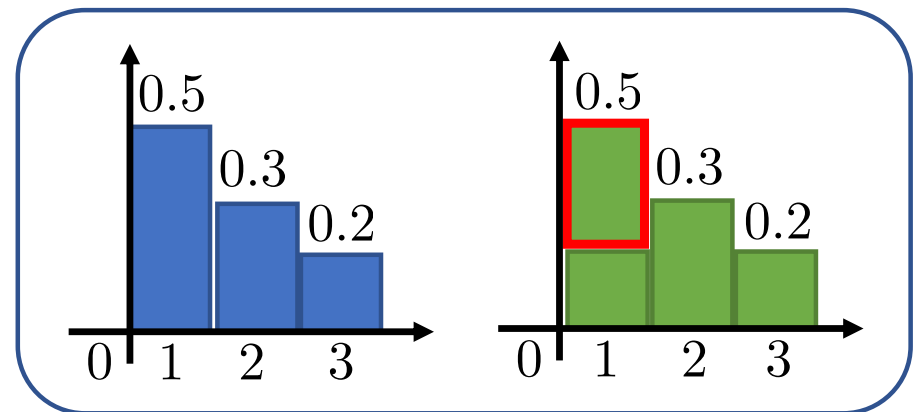
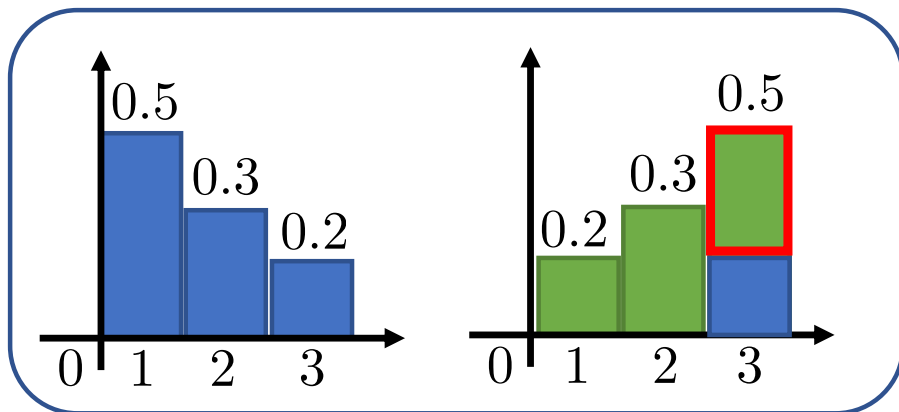
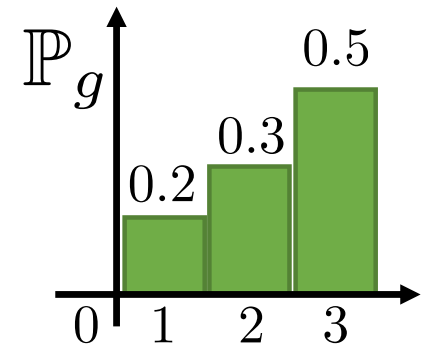
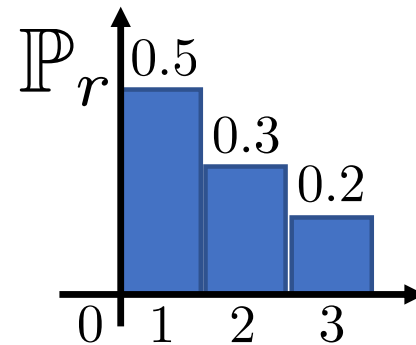
- x と y についての上記を満たす同時分布として解釈できる.
- $\gamma(x, y)$ は \mathbb{P}_g を \mathbb{P}_r に変形するために x から y へ輸送する土の量.

同時分布 γ の解釈

- 簡単のため右記の離散分布を考える.
- \mathbb{P}_g を変形させることで \mathbb{P}_r の形にする.

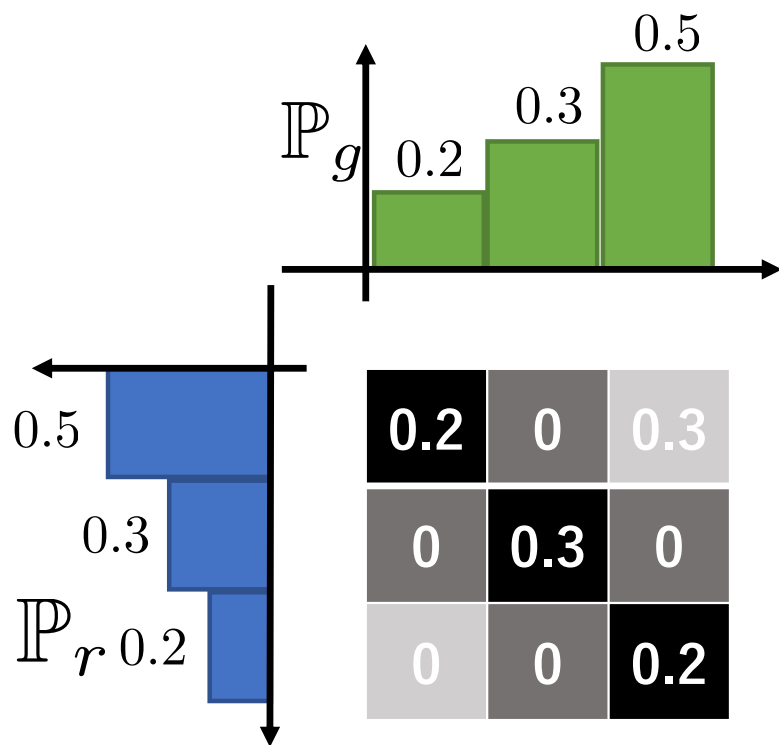
$\gamma(x, y)$ は x から y へ 輸送される土の量.

輸送方法として次の手順が考えられる.



3のビンの赤色の量を1のビンへ輸送することで形が一致する.
このとき $\gamma(x, y)$ は →

➤ 先の輸送方法 $\gamma(x, y)$ は下の 2 次元ヒストグラム(表)で表現できる.



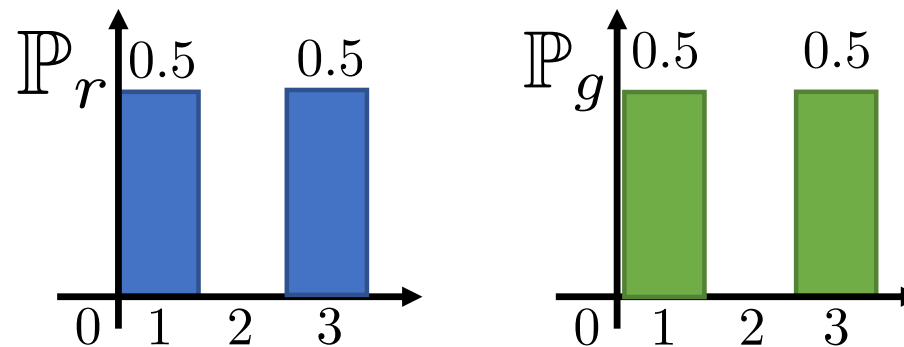
$$\sum_x \gamma(x, y) = \mathbb{P}_r(y) \quad \sum_y \gamma(x, y) = \mathbb{P}_g(x)$$

- もちろん上記の制約を満たす.
- 対角成分(x=y)は輸送しない量となる.
- 違う輸送方法については違う表になる.

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|].$$

- 対角成分に大きい数値があるほどEMDは小さくなる.

- 次の二つの同じ形の分布を考える.
 - 最適な輸送方法は”輸送しないこと”であることは自明.

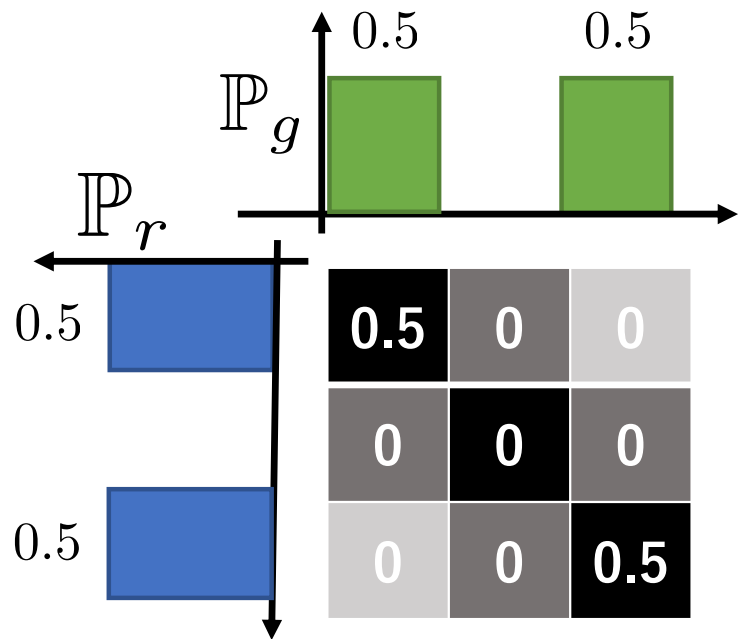


- ただ、1のビンと3のビンを入れ替えるという無駄な輸送方法を考えられる.

最適な輸送方法と無駄な輸送方法を比較してEMDがどうなるか確認する.

※(EMDは最適な輸送方法での距離なので無駄な輸送方法を用いた値はEMDではないことに注意)

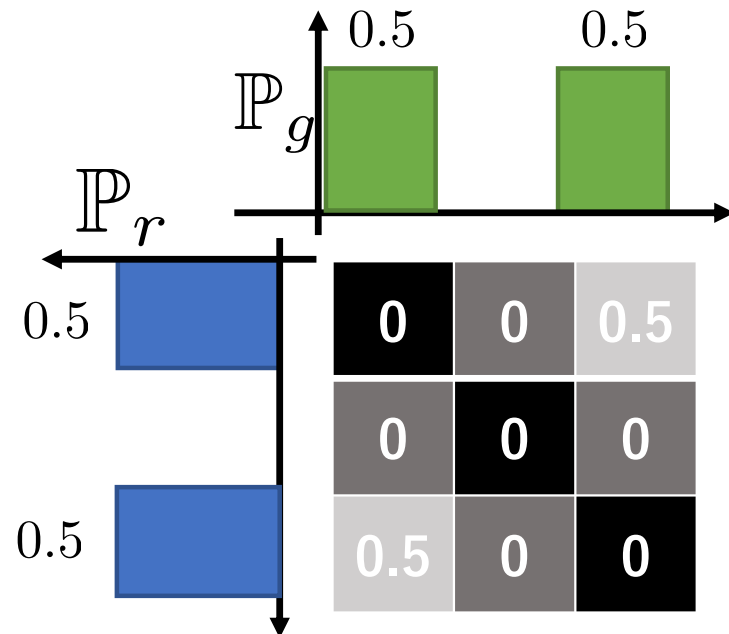
最適



$$\sum_{(x,y)} \{ \|x - y\| \gamma(x,y) \} = 0$$

こっちがEMDとなる。

無駄な輸送方法



$$\sum_{(x,y)} \{ \|x - y\| \gamma(x,y) \}$$

$$= \|3 - 1\| \cdot 0.5 + \|1 - 3\| \cdot 0.5 = 2$$

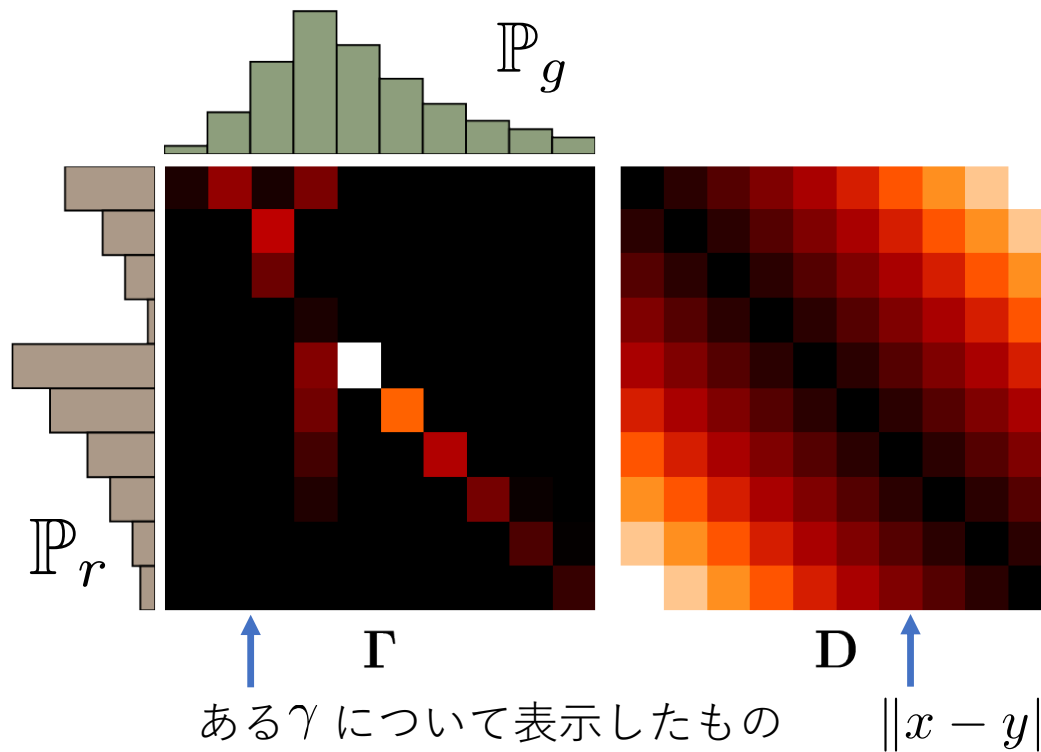
$$\text{EMD}(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi} \sum_{x, y} \|x - y\| \gamma(x, y) = \inf_{\gamma \in \Pi} \mathbb{E}_{(x, y) \sim \gamma} \|x - y\|$$

$$\text{EMD}(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi} \langle \mathbf{D}, \mathbf{\Gamma} \rangle_F$$

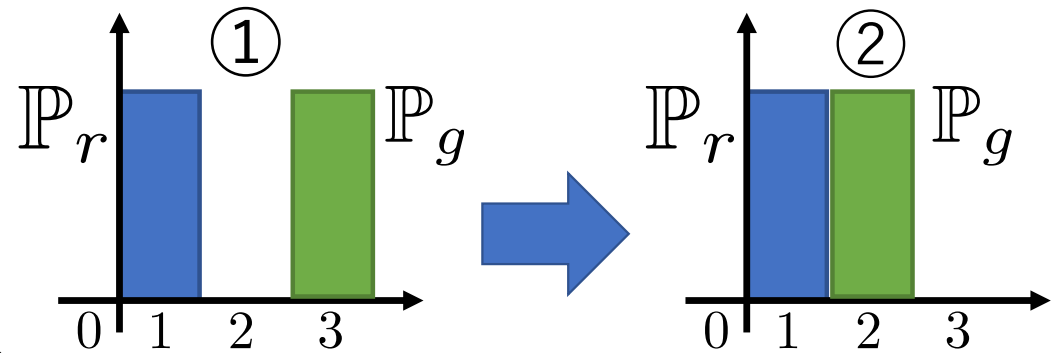
$\langle \mathbf{D}, \mathbf{\Gamma} \rangle_F$ はフロベニウス積

可能な限り対角成分に大きい数値を当てた γ が最適解(EMD)を与える γ となる.

横軸x縦軸y.



- 右のように緑の分布を近づける.
- ①②それぞれにおけるEMDは次のようになる.



$$\text{EMD}(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi} \sum_{x,y} \|x - y\| \gamma(x, y)$$

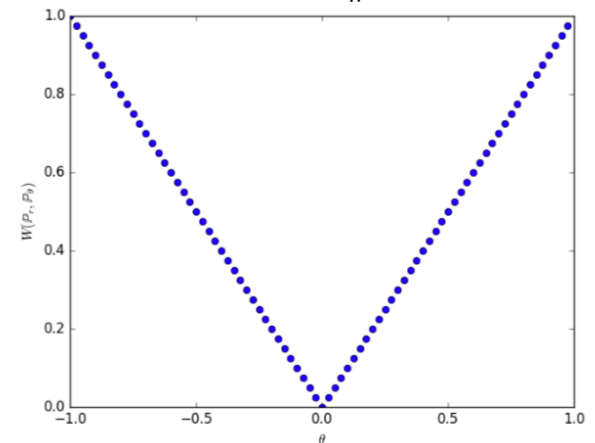
① $\text{EMD}(\mathbb{P}_r, \mathbb{P}_g) = \|3 - 1\| \cdot 1 = 2$

② $\text{EMD}(\mathbb{P}_r, \mathbb{P}_g) = \|2 - 1\| \cdot 1 = 1$

- 分布が一致したときもちろん0.

絶対値がそのまま現れて定数
にならず不連続でもない.

緑の分布を動かし続けた
ときのEMDの値



1. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems* (pp. 2672-2680).
2. Goodfellow, I. (2016). NIPS 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*.
3. Mohamed, S., & Lakshminarayanan, B. (2016). Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*.
4. Radford, A., Metz, L., & Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
5. Mirza, M., & Osindero, S. (2014). Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*.
6. Karras, T., Aila, T., Laine, S., & Lehtinen, J. (2017). Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*.
7. Arjovsky, M., Chintala, S., & Bottou, L. (2017). Wasserstein gan. *arXiv preprint arXiv:1701.07875*.