

Generative Adversarial Networks

by Paul Hand
Northeastern University

Outline

- GANs - examples and properties
- Minimax Formulation and theory
- Wasserstein GANs
- Challenges

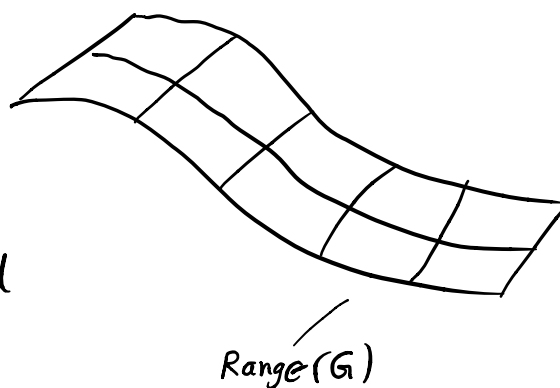
Generative Adversarial Networks (Goodfellow et al. 2014)

Generative model trained in a game-theoretic adversarial way

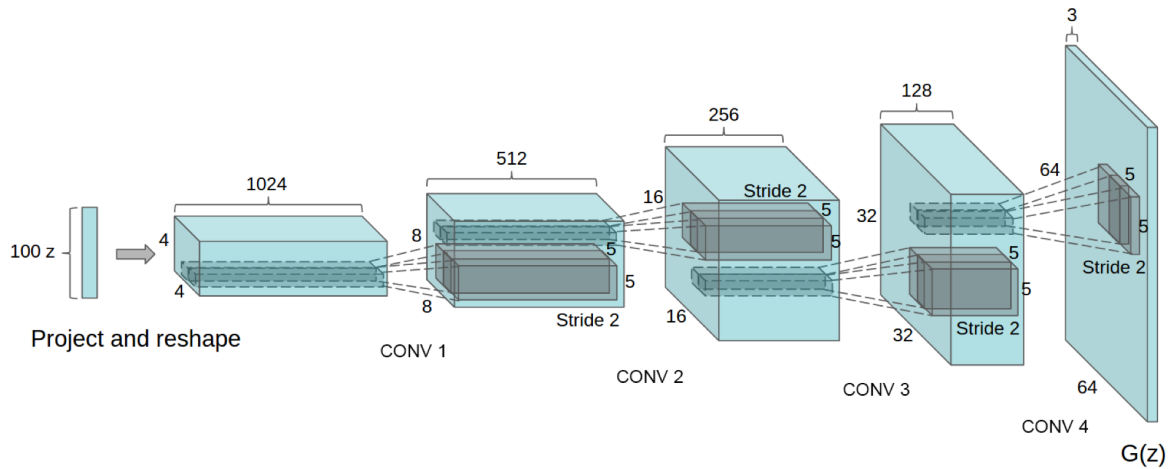
$G: \mathbb{R}^k \rightarrow \mathbb{R}^n$ st if $z \sim \mathcal{N}(0, \mathbf{I}_k)$ then $G(z)$ samples from a learned data distribution

latent space image space

While G induces a distribution on \mathbb{R}^n , we will not attempt to maximize data likelihood



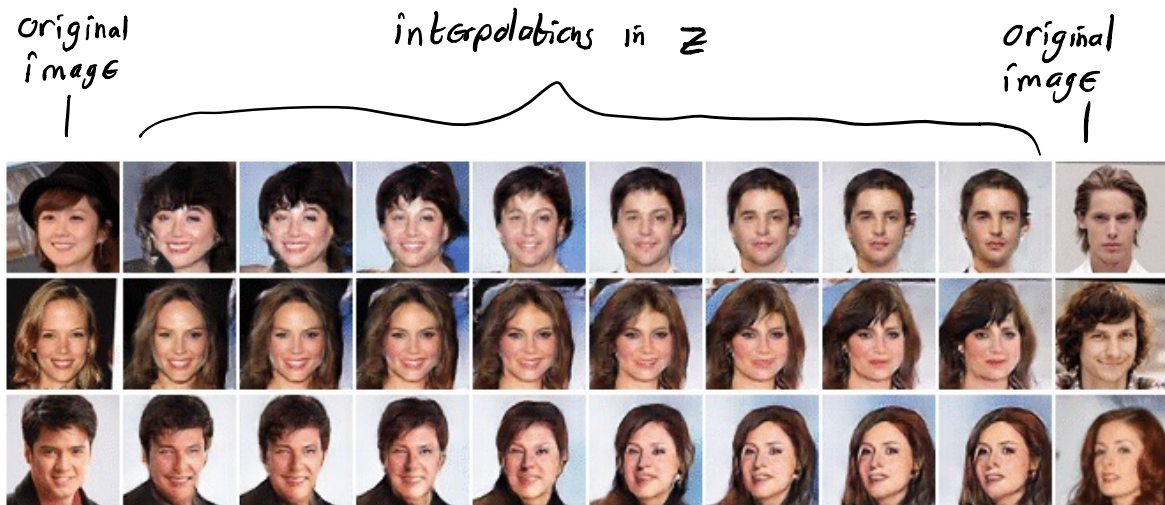
Example architecture (DCGAN) (Radford et al. 2016)



Synthetic Samples when trained on LSUN Bedrooms:

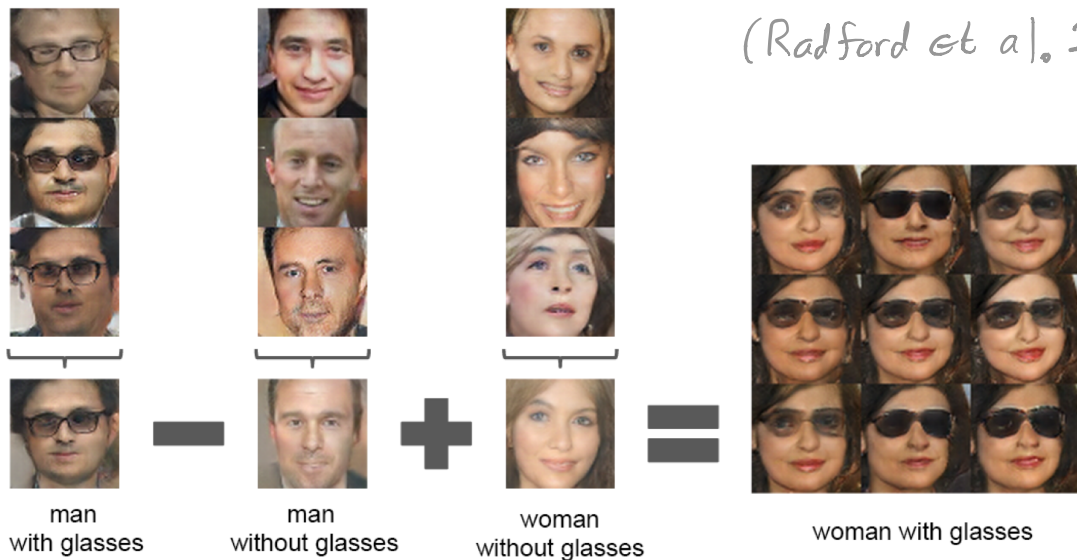


Can interpolate in latent space?



(ulyahov et al. 2017)

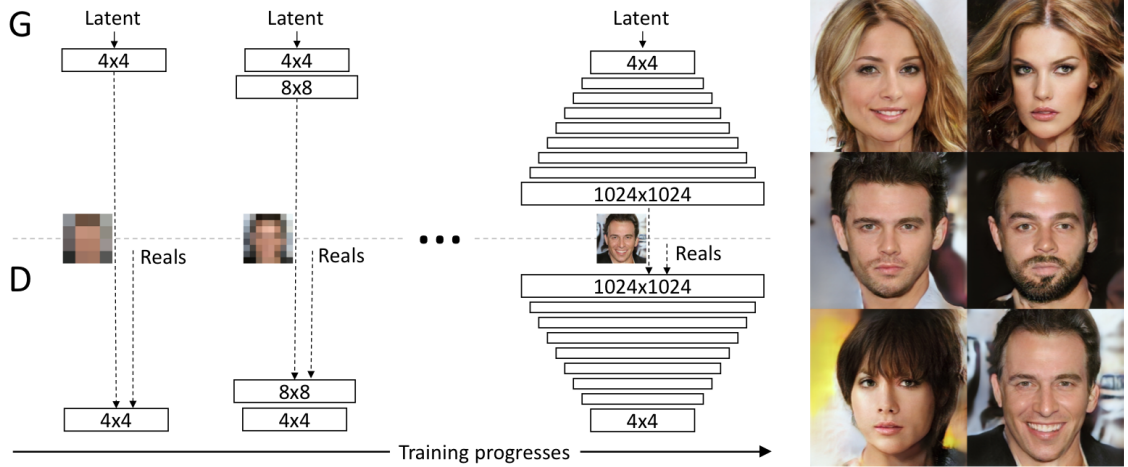
There is semantically meaningful arithmetic in latent space?



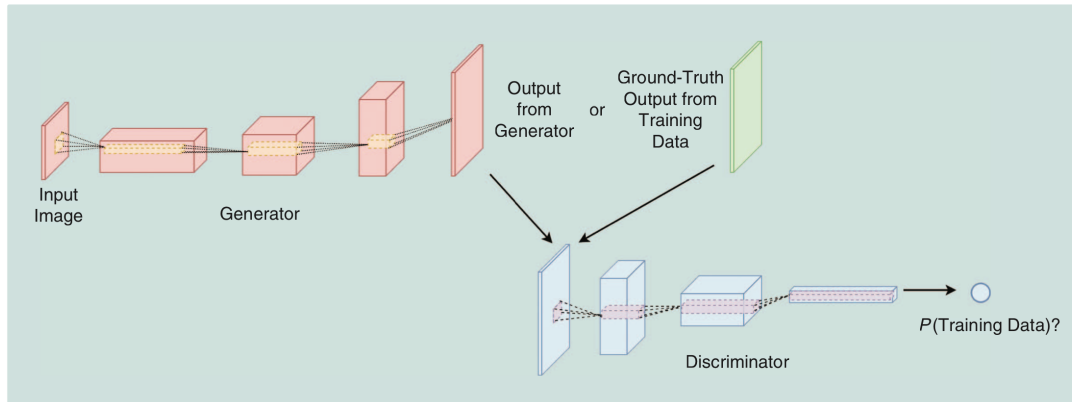
There is a direction in z corresponding to having glasses

GANs have been trained that can generate photorealistic faces

(Karras et al. 2018)



Idea: Train a model by trying to fool a concurrently trained discriminator



(Lucas et al. 2018)

Formulation of GAN training as minimax optimization

Let P_d denote data distribution
 $P_z \in \mathcal{N}(0, I_k)$

Let $G: \mathbb{R}^k \rightarrow \mathbb{R}^n$ be the generator
 $D: \mathbb{R}^n \rightarrow [0, 1]$ be $P(\text{input is real})$

Value function

$$V(D, G) = \mathbb{E}_{x \sim P_d} \log D(x) + \mathbb{E}_{z \sim P_z} \log(1 - D(G(z)))$$

Why optimize this?

it is the negative cross-entropy loss
 but label = real when $X \sim P_x$
 and label = not real when $Z \sim P_z$

Cross entropy loss

$$l_{CE}(P, q) = - \sum_{s \in S} P(s) \log q(s) = - \mathbb{E}_P(\log q)$$

$\underbrace{\hspace{1.5cm}}_{\text{r.v.s over } S}$

Minimax formulation

$$\min_G \max_D \mathbb{E}_{x \sim P_x} \log D(x) + \mathbb{E}_{z \sim P_z} \log(1 - D(G(z)))$$

\downarrow D wants to maximize neg. cross-entropy
 \downarrow G wants the opposite

Minibatch Stochastic Gradient Descent Algorithm

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:

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

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:

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

end for

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

(Goodfellow et al. 2014)

Why is the GAN value function the right thing to optimize?

Claim: For fixed G , the optimal D is

$$D_G^*(x) = \frac{P_d(x)}{P_d(x) + P_g(x)}$$

Proof:
$$V(G, D) = \mathbb{E}_{x \sim P_d} \log D(x) + \mathbb{E}_{z \sim P_z} \log(1 - D(G(z)))$$
$$= \mathbb{E}_{x \sim P_d} \log D(x) + \mathbb{E}_{x \sim P_g} \log(1 - D(x))$$

distribution
induced by generator

$$= \int_{\mathcal{X}} (P_d(x) \log D(x) + P_g(x) \log (1-D(x))) dx$$

To find max over D :

Use Variational Calculus and differentiate with respect to D and set equal to 0

$$\frac{P_d(x)}{D(x)} - \frac{P_g(x)}{1-D(x)} \equiv 0$$

$$\Rightarrow D^*(x) = \frac{P_d(x)}{P_d(x) + P_g(x)} \quad \square$$

Theorem: The global minimum of

$$C(G) = \max_D V(G, D)$$

is unique and achieved iff $P_g = P_d$.

Proof: By previous claim,

$$C(G) = \mathbb{E}_{x \sim P_d} \log D_G^*(x) + \mathbb{E}_{x \sim P_g} \log (1 - D_G^*(x))$$

$$\begin{aligned}
&= \mathbb{E}_{x \sim P_d} \log P_d \frac{2}{P_d + P_g} + \mathbb{E}_{x \sim P_g} \log P_g \frac{2}{P_d + P_g} - \log 4 \\
&= -\log 4 + D_{KL} \left(P_d \parallel \frac{P_d + P_g}{2} \right) + D_{KL} \left(P_g \parallel \frac{P_d + P_g}{2} \right) \\
&\quad \quad \quad \backslash \quad \quad \quad / \\
&\quad \quad \quad \text{nonnegative and} \\
&\quad \quad \quad 0 \text{ iff } P_d = P_g
\end{aligned}$$

Limits on this theory:

Nonparametric, infinite capacity models
(all probability distributions)

Does not assure the minimax problem
can be solved to global optimality

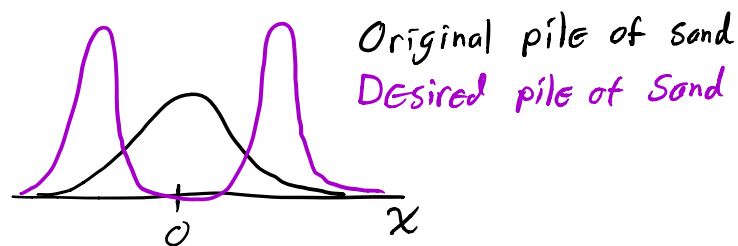
Wasserstein GAN

(Arjovsky et al. 2017)

Goal: minimize "distance" between P_d and P_g

Use Earth mover distance
(Wasserstein-1 distance)

Illustration:



Move each grain such that average distance moved is minimized

Formally,

$$W(P_d, P_g) = \inf_{\gamma \in \Pi(P_d, P_g)} \mathbb{E}_{(x,y) \sim \gamma} \|x-y\|$$

$$\text{w/ } \Pi(P_d, P_g) = \left\{ \begin{array}{l} \text{joint distributions on } (x,y) \\ \text{s.t. marginals are } P_d \text{ and } P_g \end{array} \right\}$$

Why minimize EMD?

Plain GAN (earlier) roughly minimizes

$$D_{KL}(P_d \parallel \frac{P_d + P_g}{2}) + D_{KL}(P_g \parallel \frac{P_d + P_g}{2}) = JS(P_d, P_g)$$

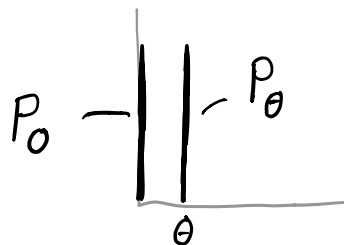
Jensen-Shannon
divergence

This is not continuous in P_d and P_g , but
EMD is.

Example:

Consider uniform distribution over
the 2d line segment

$$P_\theta = \{(x, y) \mid 0 \leq y \leq 1\} \subset \mathbb{R}^2$$



$$KL(P_0, P_\theta) = \begin{cases} \infty & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$JS(P_0, P_\theta) = \begin{cases} \log 2 & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$W(P_0, P_\theta) = |\theta|$$

As $\theta \rightarrow 0$, only $W(P_0, P_\theta) \rightarrow 0$.

Approximating EMD w/ nets

By Kantorovich-Rubinstein duality

$$W(P_d, P_g) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P_d} f(x) - \mathbb{E}_{x \sim P_g} f(x)$$

Lipschitz constant: $\|f\|_L = \sup_{x \neq y} \frac{\|f(x) - f(y)\|}{\|x - y\|}$

At the expense of a factor of K ,
can take sup over $\|f\|_L \leq K$

To estimate $W(P_d, P_g)$:

$$\max_{w \in W} \mathbb{E}_{x \sim P_d} f_w(x) - \mathbb{E}_{z \sim P_g} f_w(G_\theta(z))$$

where f_w are neural nets w/ parameters

w in a compact set W .

eg each weight is in $[-0.01, 0.01]$

WGAN formulation

$$\min_w \max_{\theta} \mathbb{E}_{x \sim p_d} f_w(x) - \mathbb{E}_{z \sim p_z} f_w(G_{\theta}(z))$$

Call f_w the "critic"

Algorithm:

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

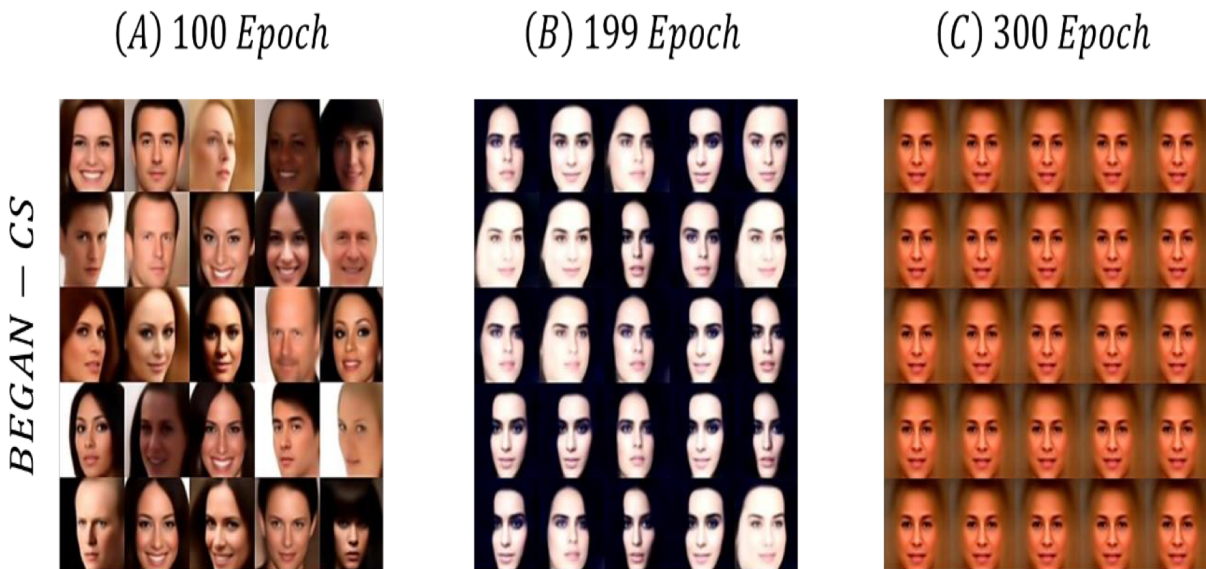
Require: α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: w_0 , initial critic parameters. θ_0 , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_{\theta}(z^{(i)}))]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_{\theta} \leftarrow -\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m f_w(g_{\theta}(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_{\theta})$ 
12: end while
```

Challenges with GANs:

- Difficulty in training (eg # D updates per G update)
- Mode collapse



(Park et al. 2020)

- No evaluation metric
- No likelihood estimates
- Difficult to invert

$$\min_z \|G(z) - y\|^2$$