

Generative adversarial networks

Victor Kitov

v.v.kitov@yandex.ru



Table of Contents

1 Simple GAN

- Model
- Practical recommendations

2 Deep convolutional GAN

3 Semi-supervised learning with GAN

4 Conditional GAN

5 CycleGAN

6 Progressive GAN

Simple GAN

Model

1 Simple GAN

- Model
- Practical recommendations

Density estimation

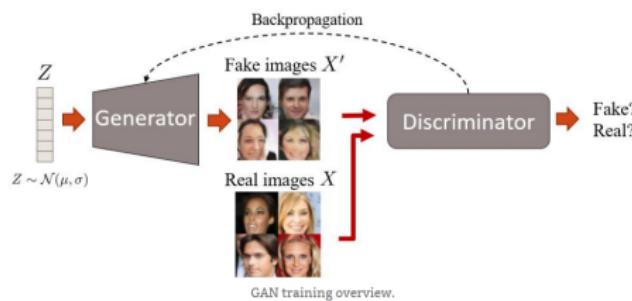
- Consider i.i.d. sample x_1, \dots, x_N , $x_i \sim p(x)$.
- Want to estimate $p(x)$. Approximate it with $q_\theta(x)$.
- **Explicit:** find θ from ML

$$\hat{\theta} = \arg \max_{\theta} \frac{1}{N} \sum_{n=1}^N \ln q_\theta(x_n) = \arg \min_{\theta} KL(p||q_\theta)$$

- **Implicit density estimation** (used in GAN):
 - do not estimate function $q_\theta(x)$ explicitly
 - instead make a sampler from $q_\theta(x)$
- We replace hard problem (estimation of multidimensional $p(x)$) with simple one (estimating a classifier).

Intuition of adversarial learning

Generative adversarial learning for images:



Analogy for bank and a money counterfeiter (having a spy in the bank).

- they compete, until money counterfeiter learns to make perfect money replicas!

Seminal paper on GAN¹

- Two networks:
 - generator $G(z) : Z \rightarrow X$
 - outputs generated object x
 - discriminator $D(x) : X \rightarrow [0, 1]$
 - probability that x **is real** rather than generated by G .

¹[Link to paper.](#)

Seminal paper on GAN¹

- Two networks:
 - generator $G(z) : Z \rightarrow X$
 - outputs generated object x
 - discriminator $D(x) : X \rightarrow [0, 1]$
 - probability that x is **real** rather than generated by G .
- Define
 - $p(x)$ - true data distribution (from training set)
 - $q(x)$ - generated data distribution $x \sim q(x) = G(z)$, $z \sim p(z)$.
 - $p(z)$ - standard distribution (Gaussian or other-spherical recommended).

¹[Link to paper.](#)

Game

- Probability that x is correctly classified by discriminator:

$$\begin{cases} D(x), & x \text{ is real} \\ 1 - D(x), & x \text{ is fake} \end{cases}$$

Game

- Probability that x is correctly classified by discriminator:

$$\begin{cases} D(x), & x \text{ is real} \\ 1 - D(x), & x \text{ is fake} \end{cases}$$

- Log-probability of correct classification by discriminator:

- given that $p(\text{real}) = p(\text{fake}) = \frac{1}{2}$

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

Game

- Probability that x is correctly classified by discriminator:

$$\begin{cases} D(x), & x \text{ is real} \\ 1 - D(x), & x \text{ is fake} \end{cases}$$

- Log-probability of correct classification by discriminator:

- given that $p(\text{real}) = p(\text{fake}) = \frac{1}{2}$

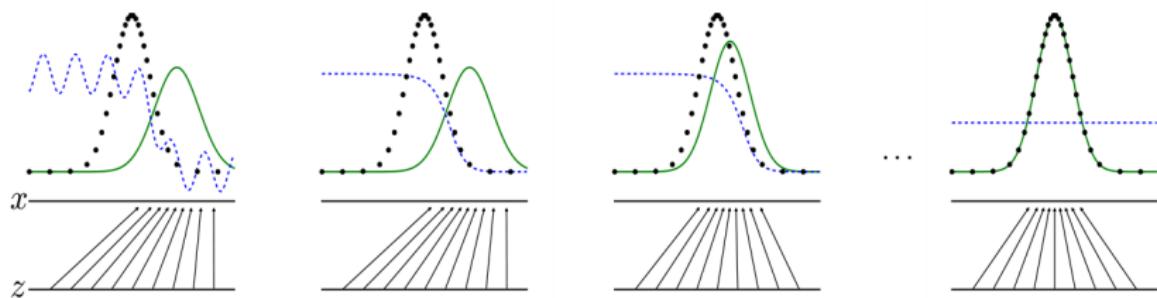
$$V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

- D and G play two-player game with minimax function $V(G, D)$:

$$\min_G \max_D V(D, G)$$

Illustration of incremental learning

Incremental learning of D and G :



black dotted: $p(x)$ - density of true samples

green: $q(x)$ - density of fake samples

blue dashed: $D(x) = p(x \text{ is true} | x)$

Losses

D task (for fixed G):

$$\mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \max_D$$

G task (for fixed D):

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \min_G$$

Losses

D task (for fixed G):

$$\mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \max_D$$

G task (for fixed D):

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \min_G$$

- Thus $G(z)$ learns to sample realistic samples from $p(x)$.
- D and G should be trained synchronously.
 - if too much training of D , G will get zero-gradient, no feedback!
 - if too much training of G , it may converge to $\arg \max_x D(x)$.

Algorithm

for number of training iterations **do**

for k steps **do**

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

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

end for

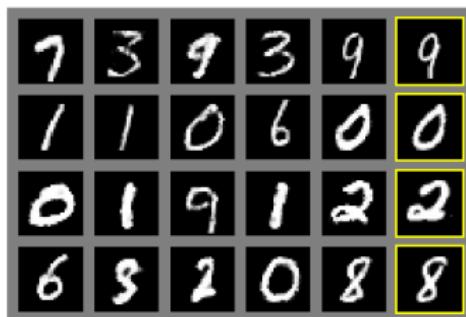
- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

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

end for

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

Generated images



a)



b)



Interpolation in latent space

Linear interpolation of objects in latent space:



Optimal G^* given fixed D

- Optimal G^* given fixed D gives delta function $\delta(x - \arg \max_x D(x))$.
 - called *mode collapse*.
- Even if didn't converge to mode, tends to collapse to regions with high $p(x)$.
 - no diversity of samples!

Optimal D^* given fixed G

Suppose D is flexible enough to take arbitrary values.

Theorem: For fixed G optimal discriminator is:

$$D^*(x|G) = \frac{p(x)}{p(x) + q(x)}$$

Optimal D^* given fixed G

Suppose D is flexible enough to take arbitrary values.

Theorem: For fixed G optimal discriminator is:

$$D^*(x|G) = \frac{p(x)}{p(x) + q(x)}$$

Proof:

$$\begin{aligned} V(G, D) &= \int_x p(x) \log(D(x)) dx + \int_z p(g(z)) \log(1 - D(g(z))) dz = \\ &= \int_x p(x) \log(D(x)) dx + q(x) \log(1 - D(x)) dx \end{aligned}$$

Applying for $\forall x \in X \arg \max_{d \in \mathbb{R}} \{p \log(d) + q \log(1 - d)\} = \frac{p}{p+q}$ for any $p, q \in \mathbb{R}$ obtain

$$\arg \max_D V(G, D) = \frac{p(x)}{p(x) + q(x)}$$

Optimal G^* given D^*

Theorem: given optimal D^* , optimal G^* should yield $q(x) = p(x)$.

Optimal G^* given D^*

Theorem: given optimal D^* , optimal G^* should yield $q(x) = p(x)$.

Proof:

Let G^* yield $q^*(x) = p(x)$. Then

$$\begin{aligned} V(G^*, D^*) &= \int_x p(x) \log \frac{p(x)}{p(x) + q^*(x)} dx + \int_x q(x) \log \frac{q(x)}{p(x) + q^*(x)} dx \\ &= \int_x p(x) \log \frac{1}{2} dx + \int_x q(x) \log \frac{1}{2} dx = 2 \log \frac{1}{2} = -2 \log 2 \end{aligned}$$

Optimal G^* given D^*

And $V(G^*, D^*) = \min_G V(G, D^*)$, because

$$\begin{aligned}
 & V(G, D^*) - V(G^*, D^*) \\
 &= \int_x p(x) \log \frac{p(x)}{p(x) + q(x)} dx + \int_x q(x) \log \frac{q(x)}{p(x) + q(x)} dx \\
 &\quad + \int_x p(x) \log 2 dx + \int_x q(x) \log 2 dx \\
 &= \int_x p(x) \log \frac{p(x)}{\frac{p(x)+q(x)}{2}} dx + \int_x q(x) \log \frac{q(x)}{\frac{p(x)+q(x)}{2}} dx \\
 &= KL \left(p(x) \parallel \frac{p(x) + q(x)}{2} \right) + KL \left(q(x) \parallel \frac{p(x) + q(x)}{2} \right) \geq 0
 \end{aligned}$$

1 Simple GAN

- Model
- Practical recommendations

Finding unstable saddle-point

D task (for fixed G):

$$\mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \max_D$$

G task (for fixed D):

$$\mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))] \rightarrow \min_G$$

Finding saddle-point $\min_G \max_D V(D, G)$ is very unstable.

- requires a lot of fine-tuning, see [link to recommendations](#).

Amplification of gradients

Problem:

- On early iterations G generates poor fakes.
- It becomes very easy for D to discriminate.
- So $D(G(z)) \approx 0 \Rightarrow \nabla_{\theta_G} \log(1 - D(G(z))) \approx 0$
 \Rightarrow no training in G .

Amplification of gradients

Problem:

- On early iterations G generates poor fakes.
- It becomes very easy for D to discriminate.
- So $D(G(z)) \approx 0 \Rightarrow \nabla_{\theta_G} \log(1 - D(G(z))) \approx 0$
 \Rightarrow no training in G .

Solution:

- improve G from another criterion:

$$\mathbb{E}_{z \sim p(z)} [\log(D(G(z)))] \rightarrow \max_G$$

- $\log(\cdot)$ exponentially amplifies small changes in gradient near zero.
 - G gets non-degenerate feedback.
 - D, G game stops being game with zero-sum.

Table of Contents

- 1 Simple GAN
- 2 Deep convolutional GAN
- 3 Semi-supervised learning with GAN
- 4 Conditional GAN
- 5 CycleGAN
- 6 Progressive GAN

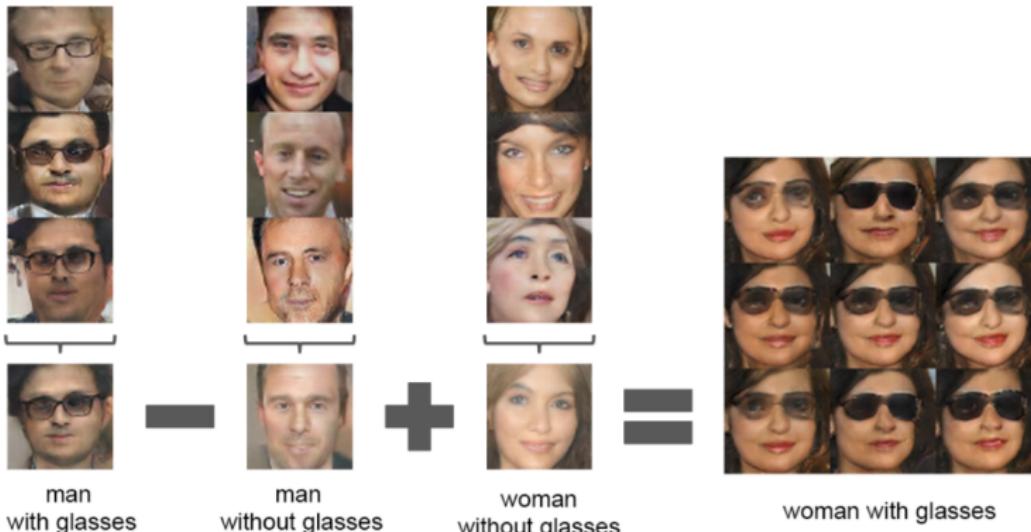
Deep convolutional GAN (DCGAN)²

Generates photorealistic 64x64 images (bedrooms in this case).



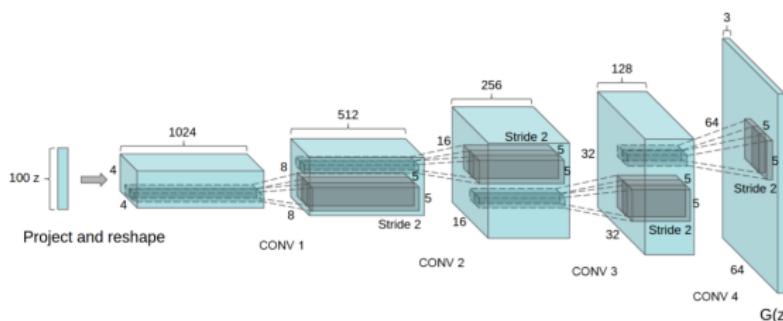
²<https://arxiv.org/pdf/1511.06434.pdf>

Latent space arithmetics



Generator

Uses fully-convolutional generator:



Architecture guidelines for stable DCGANs

- ➊ Replace any pooling layers with fractional-strided convolutions in G and strided convolutions in D .
- ➋ Use LeakyReLU activation in the discriminator for all layers.
- ➌ Use batchnorm in both the generator and the discriminator.
- ➍ Remove fully connected hidden layers for deeper architectures.
- ➎ Use ReLU activation in generator for all layers except for the output, which uses Tanh.

Intuition: 1 and 2 allow better propagation of gradients.

Table of Contents

- 1 Simple GAN
- 2 Deep convolutional GAN
- 3 Semi-supervised learning with GAN
- 4 Conditional GAN
- 5 CycleGAN
- 6 Progressive GAN

Semi-supervised learning with GAN³

- Semisupervised GAN (SGAN):
 - classifier and discriminator are united
 - classifier outputs $C + 1$ probabilities:
$$[p(y = 1|x), \dots, p(y = C|x), p(x \text{ was generated}|x)]$$
- **Weight sharing** between discriminator and classifier help each other.

³[Link to paper.](#)

Algorithm of SGAN

Algorithm 1 SGAN Training Algorithm

Input: I : number of total iterations

for $i = 1$ **to** I **do**

 Draw m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

 Draw m examples $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ from data generating distribution $p_d(x)$.

 Perform gradient descent on the parameters of D w.r.t. the NLL of D/C's outputs on the combined minibatch of size $2m$.

 Draw m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.

 Perform gradient descent on the parameters of G w.r.t. the NLL of D/C's outputs on the minibatch of size m .

end for

SGAN experiments

SGAN converges faster:

Generated MNIST images by SGAN (left) and GAN (right) after 2 MNIST epochs



SGAN experiments

Semi-supervised learning improves accuracy for small training sets.

Accuracy comparisons of supervised and semisupervised classifier on MNIST:

EXAMPLES	CNN	SGAN
1000	0.965	0.964
100	0.895	0.928
50	0.859	0.883
25	0.750	0.802

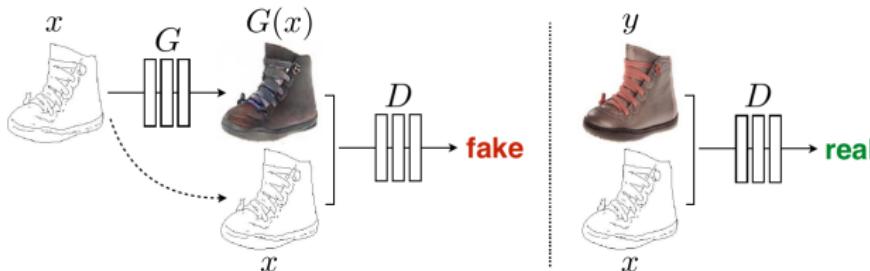
Table of Contents

- 1 Simple GAN
- 2 Deep convolutional GAN
- 3 Semi-supervised learning with GAN
- 4 Conditional GAN
- 5 CycleGAN
- 6 Progressive GAN

Conditional GAN⁴

- z - random state from standard distribution $p(z)$
- x - source version of object.
- y - target version of object
- $\hat{y} = G(x, z)$ - reconstruction of target

Conditional GAN:



⁴ Isola et al. 2017.

Loss functions

ConditionalGAN (cGAN) objective:

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,y} [\log D(x, y)] + \mathbb{E}_{x,z} [\log (1 - D(x, G(x, z)))]$$

- score for D , loss for G
- promotes realism of generated \hat{y}
- power of GAN: realism loss is determined automatically, not hard-coded.

Reconstruction loss for D :

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z} \|y - G(x, z)\|_1$$

Generator solves:

$$G^* = \arg \min_G \max_D \{\mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)\}$$

Results for different X, Y

Labels to Street Scene

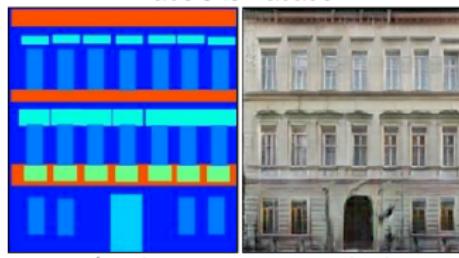


Aerial to Map

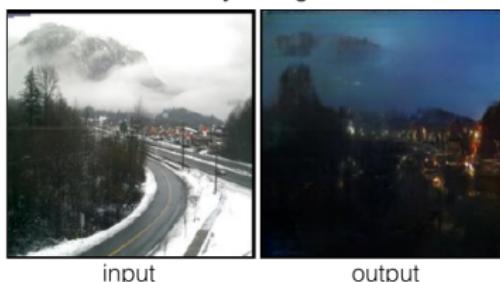
The figure consists of two side-by-side maps. The left map is a detailed satellite view of a residential neighborhood, showing individual houses with white roofs, green lawns, and a complex network of streets. The right map is a simplified version of the same area, where the buildings and green spaces are uniformly colored in a light gray shade, while the streets remain white lines. This visual comparison highlights how different levels of detail can be extracted from a single dataset.

Results for different X, Y

Labels to Facade



Day to Night



Results for different X, Y

BW to Color



input

output

Edges to Photo

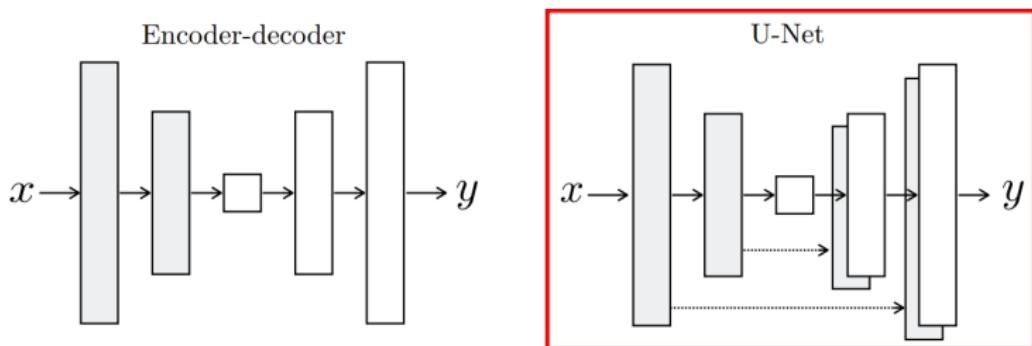


input

output

Comments

- Works for any X, Y .
 - in particular, may switch input and output spaces.
- For labeled X system generalizes to human input: [online demo](#).
- Generator uses U-net architecture:



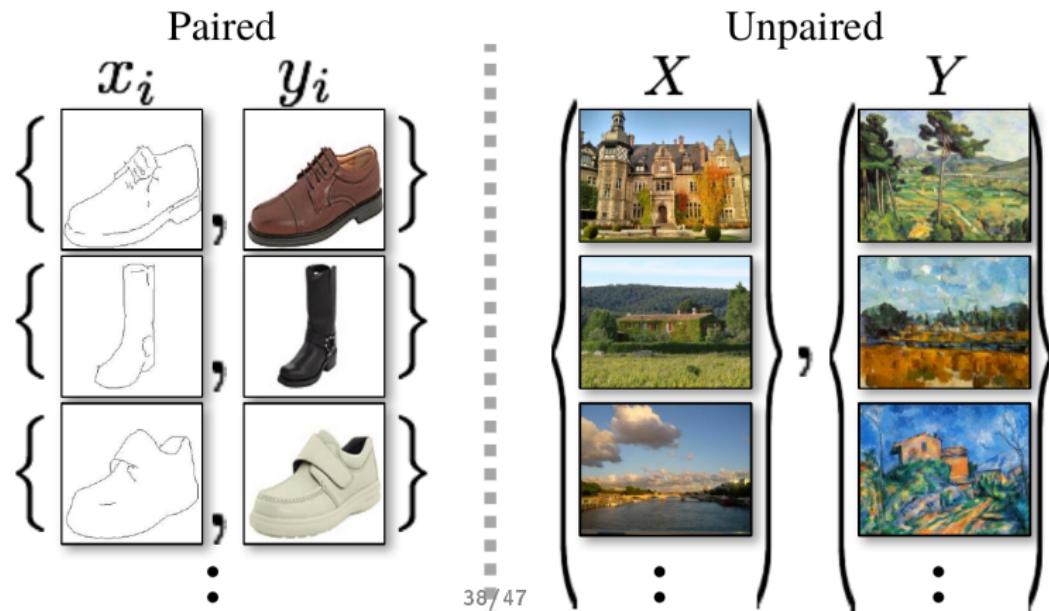
- Discriminator is applied to overlapping patches of size $N \times N$
 - allows application to arbitrary sized images.
 - $N = 70$ worked best

Table of Contents

- 1 Simple GAN
- 2 Deep convolutional GAN
- 3 Semi-supervised learning with GAN
- 4 Conditional GAN
- 5 CycleGAN
- 6 Progressive GAN

CycleGAN⁵

- 2 domains: X, Y
- **Training set is unpaired:** $x_1, x_2, \dots, x_N \in X, y_1, y_2, \dots, y_M \in Y$.
 - in contrast, conditional GAN was trained on pairs $\{(x_n, y_n)\}_{n=1}^N$



Losses

$$\begin{aligned}\mathcal{L}_{GAN}(G_{XY}, D_Y) = & \mathbb{E}_{y \sim p(y)} [\log(D_Y(y))] \\ & + \mathbb{E}_{x \sim p(x)} [\log(1 - D_Y(G_{XY}(x)))]\end{aligned}$$

$$\begin{aligned}\mathcal{L}_{GAN}(G_{YX}, D_X) = & \mathbb{E}_{x \sim p(x)} [\log(D_X(x))] \\ & + \mathbb{E}_{y \sim p(y)} [\log(1 - D_X(G_{YX}(y)))]\end{aligned}$$

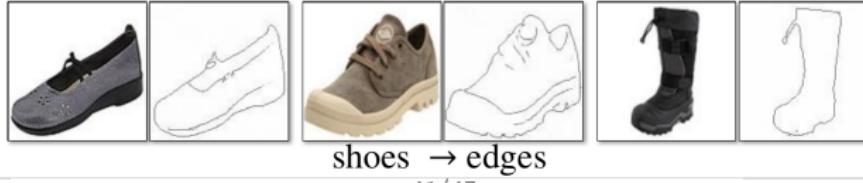
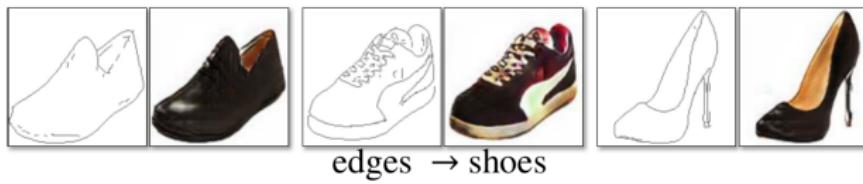
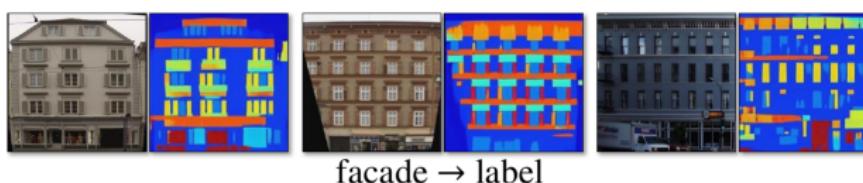
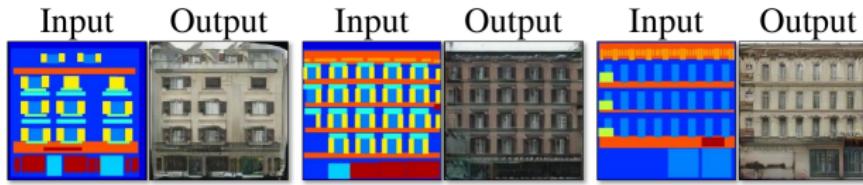
$$\begin{aligned}\mathcal{L}_{cyc}(G_{XY}, G_{YX}) = & \mathbb{E}_{x \sim p(x)} \|G_{YX}(G_{XY}(x)) - x\|_1 \\ & + \mathbb{E}_{y \sim p(y)} \|G_{XY}(G_{YX}(y)) - y\|_1\end{aligned}$$

CycleGAN

$$\begin{aligned}\mathcal{L}(G_{XY}, G_{YX}, D_X, D_Y) = & \mathcal{L}_{GAN}(G_{XY}, D_Y) + \mathcal{L}_{GAN}(G_{YX}, D_X) \\ & + \mathcal{L}_{cyc}(G_{XY}, G_{YX})\end{aligned}$$

$$G_{XY}^*, G_{YX}^* = \arg \min_{G_{XY}, G_{YX}} \max_{D_X, D_Y} \mathcal{L}(G_{XY}, G_{YX}, D_X, D_Y)$$

Demo results



Identity regularizer

- Consider paintings->photo reconstruction.
- To remove unnecessary colour change in output identity loss was added:
 - intuition: do nothing, when input is from target set.

$$\mathcal{L}_{identity} = \mathbb{E}_{x \sim p(x)} \|G_{YX}(x) - x\|_1 + \mathbb{E}_{y \sim p(y)} \|G_{XY}(y) - y\|_1$$

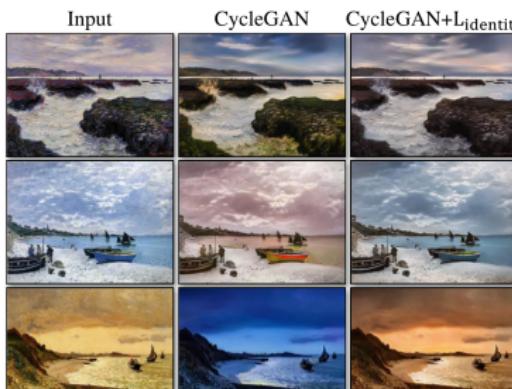


Table of Contents

- 1 Simple GAN
- 2 Deep convolutional GAN
- 3 Semi-supervised learning with GAN
- 4 Conditional GAN
- 5 CycleGAN
- 6 Progressive GAN

Progressive GAN⁶

- Technique to generate high dimensional output.

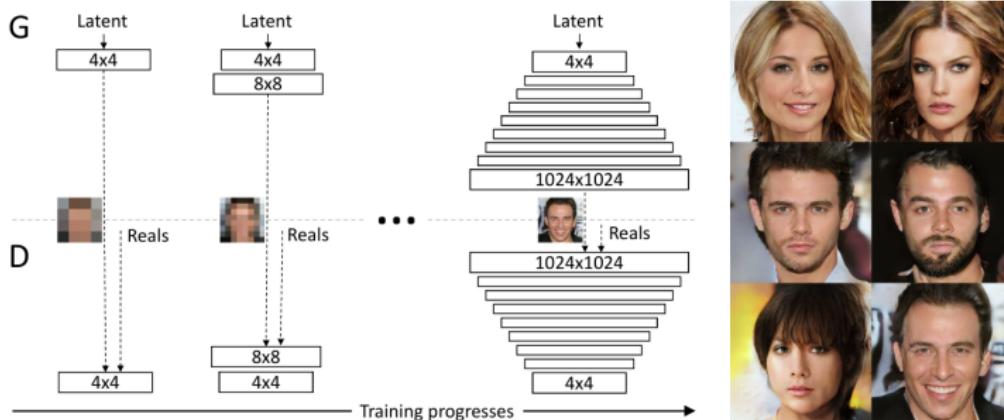
Generated photographs, using celebrities dataset



⁶Karras et al. 2017.

Training

Training is done sequentially. G, D trained on low-resolution layers, then higher resolution layers added:



Training

When new layer added, 2 branches appear (like in ResNet):

- rescale+**identity** with weight $1 - \alpha$
- rescale+**conv+nonlinearity** with weight α

During training α gradually changes from 0 to 1 (drop identity branch).

This way we gradually introduce new layer to the model.

Conclusion

- GANs - method to generate new realistic samples based on x_1, \dots, x_N .
 - requires tuning to synchronize D and G training
 - suffers from mode collapse
- DCGAN - GAN designed for images generation.
- Conditional GAN generates \hat{y} conditional on input
 - inpainting, colorization, deblurring, noise removal, etc.
- CycleGAN learns reconstructions on unpaired datasets.
- ProgressiveGAN generates high resolution images.