

Adversarial Learning

Kacper

February 21, 2016

The aim is to learn distribution p . We want to learn a generative model G , such that for a random variable Z , $G(Z) \sim p$. Or, $p = p_G$, where p_G is distribution of $G(Z)$.

Generative Adversarial Networks

artificial labels

- Random variable X , distributed according to distribution p .
- Some noise random variable Z .
- A generative model G , which transforms Z to $X' = G(Z)$, in other words $X' \sim p_G$.

Consider an artificial label $Y \in \{-1, 1\}$, distributed according to Bernoulli distribution. Consider a random vector W , such that $W_1 = Y$.

if $Y = 1$ then $W_2 = X \sim p$.

if $Y = -1$ then $W_2 = X' \sim p_G$.

I probably could have written $W = (1, X) | (-1, X')$.

Let D be discriminative function which is suppose to figure out if W_2 comes from p or p_G (if $W_2 = X'$ or $W_2 = X$).

For a random variable V , $D(W_2)$ a probability that $W_1 = 1$ (label is one).

Learning density

Consider an expression

$$E(\log(D(W_2))|Y=1) + E(\log(1-D(W_2))|Y=-1).$$

It can be written using $X, G(Z)$, as a function of D, G

$$V(D, G) = E \log(D(X)) + E \log(1 - D(G(Z)))$$

For each generator G there exists an optimal discriminator D that maximizes the above function (we will show).

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

We are going to show that G that minimizes $C(G)$ is the best generator, i.e. $p_G = p$.

Optimal discriminator

For each generator G there exists an optimal discriminator D that maximizes the above function (as we will show).

$$V(D, G) = E [\log D(X) + \log(1 - D(G(Z)))] \quad (1)$$

$$= E [\log D(X) + \log(1 - D(X'))] \quad (2)$$

$$= \int p(x) \log D(x) + p_G(x) \log(1 - D(x)) \quad (3)$$

which, for a fixed G , $V(D, G)$ reaches maximum at $D = \frac{p}{p+p_g}$

All the experimental results will be discussed at the end.

For a fixed, optimal $D = \frac{p}{p+p_g}$, V is a divergence

$$V(D, G) = E (\log D(X) + \log(1 - D(X')))) \quad (4)$$

$$E \left(\log \frac{p(X)}{p(X) + p_g(X)} + \log \frac{p_g(X')}{p(X') + p_g(X')} \right) = \quad (5)$$

$$D_{KL}(p || \frac{p + p_g}{2}) + D_{KL}(p_g || \frac{p + p_g}{2}) \quad (6)$$

Training generative neural networks via Maximum Mean Discrepancy optimization

The discriminator presented in the previous paper was $C(G) = \sup_D V(D, G)$. This paper suggests to put

$$C(G) = MMD(G(Z), X)$$

and the rest remains the same, we solve

$$\arg \min C(G) \tag{7}$$

Suppose θ is parameter of the network G . The gradient of MMD w.r.t. to θ on the smallest minibatch batch $y_1 = G(z_1), y_2 = G(z_2), x_1, x_2$, is

$$\begin{aligned} \nabla MMD \sim & \frac{\partial k(y_1, y_2)}{\partial(y_1, y_2)} \left(\frac{\partial y_1}{\partial \theta}, \frac{\partial y_2}{\partial \theta} \right) + \\ & - \frac{\partial k(x_1, y_1)}{2 \partial y_1} \frac{\partial y_1}{\partial \theta} - \frac{\partial k(x_2, y_1)}{2 \partial y_1} \frac{\partial y_1}{\partial \theta} + \\ & - \frac{\partial k(x_1, y_2)}{2 \partial y_1} \frac{\partial y_1}{\partial \theta} - \frac{\partial k(x_2, y_2)}{2 \partial y_1} \frac{\partial y_1}{\partial \theta} \end{aligned}$$

Algorithm 1 Stochastic gradient descent for MMD nets.

Initialize M, θ, α, k

Randomly divide training set X into N_{mini} mini batches

for $i \leftarrow 1, \text{number-of-iterations}$ **do**

 Regenerate noise inputs $\{w_i\}_{i=1,\dots,M}$ every r iterations

for $n_{\text{mini}} \leftarrow 1, N_{\text{mini}}$ **do**

for $m \leftarrow 1, M$ **do**

$y_m \leftarrow G_{\theta}(w_m)$

end for

 compute the n 'th minibatch's gradient

 update learning rate α (e.g., RMSPROP)

$\theta \leftarrow \theta - \alpha \nabla C_n$

end for

end for

Arthur noticed that variance of MMD should be taken into account in the objective i.e.

$$\frac{\overline{MMD}}{\sqrt{\text{var}(\overline{MMD})}}.$$

Are random variables suitable model for images? 'It is well known that, for any distribution p and any continuous distribution p_z on sufficiently regular spaces (...), there is a function G , such that $G(Z)$ '

Deep Mean Maps, digression

Not exactly what I've expected

Use random features as a layer in a network, L^i is a layer for i -th picture

$$L^i \in R^{H,W}$$

H, W is number of (super) pixels. For frequencies $\omega_1, \dots, \omega_d, \dots, \omega_D$, use feature in network

$$\mu_d^i = \sum_{w=1, h=1}^{H,W} \cos(\omega_d L_{h,w}^i + b))$$

Deep Generative Image Models using a Laplacian Pyramid of Adver- sarial Networks

The conditional generative adversarial net

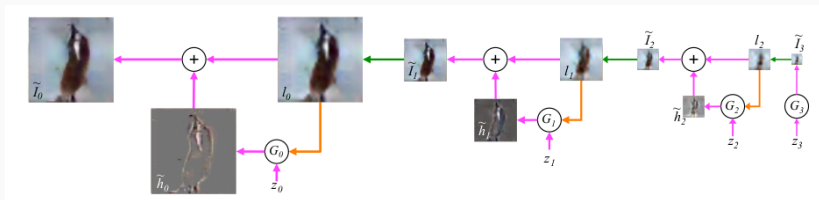
D receives additional information L as input. This might contain, say, information about the class of the training example X .

$$\arg \min_G \sup_D E \log(D(X, L)) + E \log(1 - D(G(Z), L))$$

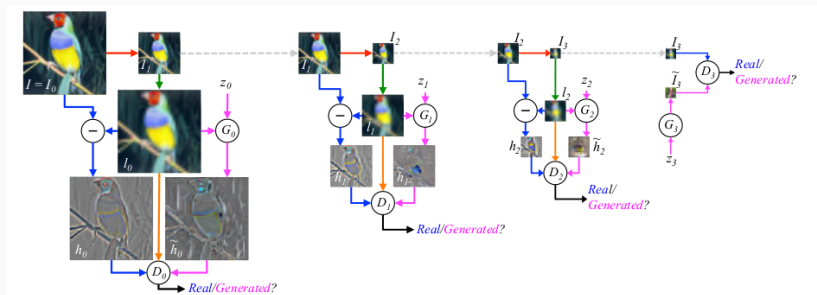
M. Mirza and S. Osindero. Conditional generative adversarial nets.

d is downsampling operation which blurs and decimates a $j \times j$ image I to $d(I)$, which is $j/2 \times j/2$. u is upsampling operation that which smooths and expands. $I_k = d^k(I)$ h_k are residuals i.e $h_k = I_k - u(d(I_k))$. They aim to learn residuals.

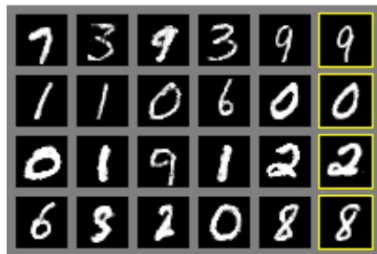
Sampling



Learning



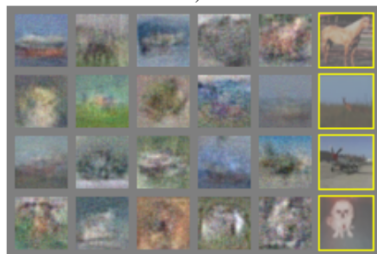
Results



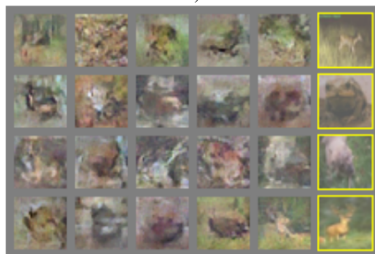
a)



b)



c)



d)

a) MNIST b) TFD c) CIFAR-10 (fully connected model) d)

9	1	6	6	5	9	3	0	3	0
6	9	1	0	7	7	6	3	7	6
9	2	8	3	8	3	8	9	8	4
9	8	0	0	3	8	5	2	6	6
6	6	7	8	6	5	0	4	6	2
8	0	0	2	3	6	2	7	0	5
8	0	1	9	5	8	0	8	2	9
6	1	4	1	2	4	3	1	5	1
4	2	6	7	6	0	7	8	4	8
5	9	7	9	5	6	0	9	6	4

1	6	0	8	1	0	8	4	8	4
7	5	3	5	0	6	5	1	8	1
6	2	9	0	9	5	9	7	9	3
8	1	7	7	8	2	1	0	5	8
2	0	5	4	3	0	1	9	1	2
5	8	1	6	7	5	8	2	0	9
1	5	2	7	9	7	6	3	3	5
1	0	0	4	5	8	0	8	1	
3	2	8	1	0	1	3	1	5	2
2	7	5	9	4	0	2	7	0	7

2	7	4	1	8	3	6	7	1	3
3	4	2	1	9	1	3	4	0	8
4	0	0	5	0	3	1	4	8	0
1	3	1	5	1	9	7	4	3	0
3	7	7	6	6	0	0	6	9	5
7	9	4	4	5	5	6	0	0	8
4	3	3	0	5	3	7	0	0	1
1	9	8	5	4	6	0	8	9	3
7	0	7	8	6	8	3	2	6	1
2	9	9	9	4	5	1	3	4	6

MNIST



CIFAR 10

LAPGAN



LSUN