

GENERATIVE ADVERSARIAL NETS

BASED ON THE ORIGINAL WORK OF GOODFELLOW ET AL.

FRANCESCO CAPORALI

Department of Mathematics, University of Turin

March 11, 2024

The authors proposed a deep generative framework developed on a game theoretic idea.

Objective: overcome the intractability issues of deep learning models that arises from conventional statistical approaches (e.g. maximum likelihood estimation).

Strategy:

	Discriminator D	Generator G
input	z : noise x : sample data	z : noise
goal	let $D(G(z))$ close to 0, let $D(x)$ close to 1	generate $G(z)$, let $D(G(z))$ close to 1

Let D, G be two maps generated through neural nets:

Discriminator. $D(x) := D(x, \theta_D)$, $D : X \rightarrow [0, 1]$.

Generator. $G(z) := G(z, \theta_G)$, $G : Z \rightarrow X$;

Given $x \in X$ (data space), $z \in Z$ (noise space), we define:

- $p_{\text{data}}(x)$, density of the input data;
 - $p_G(x)$, density of the Generator;
 - $p_{\text{noise}}(z)$, density of a prior noise.
-

D and G play a minimax game with value function $V(D, G)$,

$$V(D, G) := \mathbb{E}_{x \sim p_{\text{data}}}[\log(D(x))] + \mathbb{E}_{z \sim p_{\text{noise}}}[\log(1 - D(G(z)))].$$

We aim to find D^*, G^* s.t. $V(D^*, G^*) = \min_G \max_D V(D, G)$.

Discriminator goal

$$\max_D \underbrace{\mathbb{E}_{x \sim p_{\text{data}}} [\log(D(x))]}_{D(x) \approx 1} + \underbrace{\mathbb{E}_{z \sim p_{\text{noise}}} [\log(1 - D(G(z)))]}_{D(G(z)) \approx 0}.$$

Generator goal

$$\min_G \max_D \underbrace{\cancel{\mathbb{E}_{x \sim p_{\text{data}}} [\log(D(x))]} + \mathbb{E}_{z \sim p_{\text{noise}}} [\log(1 - D(G(z)))]}_{D(G(z)) \approx 1}.$$

Remark

Under the (unrealistic) assumption of being able to find the optimal discriminator $D_G^* := \max_D V(D, G)$, it is possible to show that iteratively alternating computations of D_G^* and updates of G (to minimize $V(D_G^*, G)$), one would obtain $p_G = p_{\text{data}}$.

Proposition 1

[Goodfellow et al.]

For a fixed generator G , the associated optimal discriminator D_G^* is

$$D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)}, \forall x \in X.$$

Proof.

- 1 Thanks to the Radon-Nikodym Theorem we have

$$\mathbb{E}_{z \sim p_{\text{noise}}}[\log(1 - D(G(z)))] = \mathbb{E}_{x \sim p_G}[\log(1 - D(x))].$$

- 2 Hence we can rewrite $V(D, G)$ as

$$V(D, G) = \int_X \log(D(x)) p_{\text{data}}(x) + \log(1 - D(x)) p_G(x) dx,$$

with $\text{supp}(p_G) = G(\text{supp}(p_{\text{noise}})) \subset X$, $\text{supp}(p_{\text{data}}) \subset X$.

- 3 Consider $f : [0, 1] \rightarrow \mathbb{R}$, $f(y) := \log(y) a + \log(1 - y) b$, with $a, b > 0$, then f is a concave map and the only maximum is realized by $y = \frac{a}{a + b}$.

- 4 Defining $\tilde{D} := \frac{p_{\text{data}}}{p_{\text{data}} + p_g}$, we have, by monotonicity of the integral, $\forall D$,

$$\begin{aligned} V(D, G) &= \int_X \log(D(x)) p_{\text{data}}(x) + \log(1 - D(x)) p_G(x) dx \leq \\ &\leq \int_X \log(\tilde{D}(x)) p_{\text{data}}(x) + \log(1 - \tilde{D}(x)) p_G(x) dx = V(\tilde{D}, G), \\ \implies \tilde{D} &= \arg \max_D V(D, G) = D_G^*. \end{aligned}$$



Theorem 1

[Goodfellow et al.]

The global minimum of the virtual training criterion given by $C(G) := \max_D V(D, G)$, is achieved if and only if $p_G = p_{\text{data}}$.

For such minimum G , $C(G) = -\log(4)$.

Proof.

- 1 Let us first rewrite the criterion in a more suitable way:

$$\begin{aligned} C(G) &= V(\arg \max_D V(D, G), G) = V(D_G^*, G) = \\ &= \mathbb{E}_{x \sim p_{\text{data}}} [\log(D_G^*(x))] + \mathbb{E}_{z \sim p_G} [\log(1 - D_G^*(x))] = \\ &= \mathbb{E}_{x \sim p_{\text{data}}} \left[\log \left(\frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_G(x)} \right) \right] + \\ &\quad + \mathbb{E}_{z \sim p_G} \left[\log \left(\frac{p_G(x)}{p_{\text{data}}(x) + p_G(x)} \right) \right]. \end{aligned}$$

- 2 Let \tilde{G} be a generator such that $p_{\tilde{G}} = p_{\text{data}}$, then $C(\tilde{G}) = 2 \log(1/2) = -\log(4)$.
- 3 Moreover, adding and subtracting $-\log(4)$, we have that

$$C(G) = -\log(4) + 2 \text{JSD}(p_{\text{data}} \| p_G), \forall G.$$

- 4 We recall the properties of the Jensen-Shannon divergence: $\forall p, q$ densities $\text{JSD}(p \| q) \geq 0$ and $\text{JSD}(p \| q) = 0$ if and only if $p = q$.
- 5 Those properties are sufficient to conclude that

$$\min_G C(G) = -\log(4) \text{ and } C(G) = -\log(4) \iff p_{\text{data}} = p_G.$$



Proposition 2

[Goodfellow et al.]

If G and D have enough capacity, and at each iteration of the algorithm, D reaches its optimum given G , and p_G is updated so as to improve the criterion

$$\mathbb{E}_{x \sim p_{\text{data}}} [\log (D_G^*(x))] + \mathbb{E}_{z \sim p_G} [\log (1 - D_G^*(x))]$$

then p_G converges to p_{data} .

Sketch.

- 1 $f(G) := \max_D V(D, G)$, $f_{D_{G_{\text{old}}}^*}(G) := V(D_{G_{\text{old}}}^*, G)$, with G_{old} which is the resulting generator after its last update.
- 2 We would like to update G descending the stochastic gradient of f because this would mean that p_G is closer to p_{data} than $p_{G_{\text{old}}}$ was.

- 3 However this is equivalent to update G descending the stochastic gradient of $f_{D_{G_{\text{old}}}^*}$ (which is what we do).

This assertion follows from the fact that the subgradients of $f_{D_{G_{\text{old}}}^*}$ are a subset of the subgradients of f : $\partial f_{D_{G_{\text{old}}}^*} \subset \partial f$ (from the convexity of $V(D, G)$).

Remark

Note that the algorithm used to train the GANs does not optimize the underlying distribution of G , p_G , but the parameters of a neural net, θ_G .

This raises theoretical problems, anyway the performances of the model are impressive, therefore GANs are reasonable to be used in practice.

Algorithm. Minibatch SGD for training D and G

```
1. function train( $D, G, data, k, n\_epochs$ )
2.   for epoch in  $1, \dots, n\_epochs$  do
3.     for  $\_$  in  $1, \dots, k$  do
4.       sample  $\mathbf{z} = (z_i)_{i=1}^m, z_i \sim p_{\text{noise}}$ 
5.       sample  $\mathbf{x} = (x_i)_{i=1}^m, x_i \sim p_{\text{data}}$ 
6.        $\theta_D = \theta_D + \lambda \nabla_{\theta_D} \hat{V}_m(D, G)$ 
7.       sample  $\mathbf{z} = (z_i)_{i=1}^m, z_i \sim p_{\text{noise}}$ 
8.        $\theta_G = \theta_G - \lambda \nabla_{\theta_G} \hat{V}_m(D, G)$ 
9.   return  $D, G$ 
```

Where the sample mean of the loss $V(D, G)$ is computed on the m -dimensional minibatches \mathbf{z}, \mathbf{x} :

$$\hat{V}_m(D, G) := \frac{1}{m} \sum_{i=1}^m \log(D(x_i)) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(G(z_i))).$$

Remark

Actual implementation:

- 6th line. It is performed one step of the ascent of the stochastic gradient via backpropagation, with loss

$$\hat{V}_m(D, G) = \frac{1}{m} \sum_{i=1}^m \log(D(x_i)) + \log(1 - D(G(z_i))) ;$$

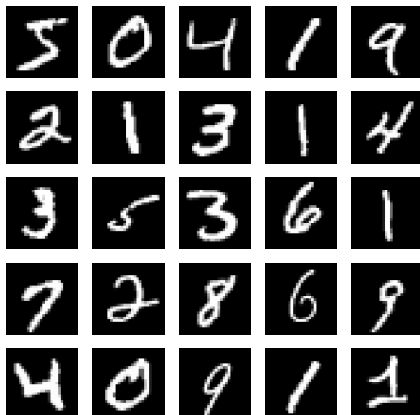
- 8th line. Instead of one step of the ascent of $\hat{V}_m(D, G)$ we, equivalently, do a step of descent with loss

$$\frac{1}{m} \sum_{i=1}^m \log(D(G(z_i))) ;$$

Note. The classical SGD is often substituted by SGD with momentum or adaptative algorithm (like ADAM).

We performed practical simulations, training and generating data using GANs (based on the work of Vandegar).

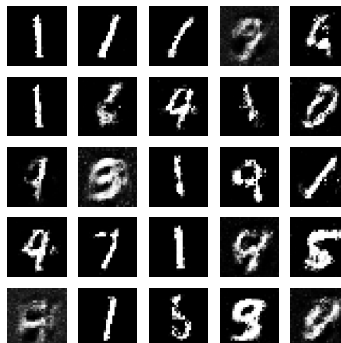
The dataset used as train set is MNIST, a database of 60000 images of handwritten digits.



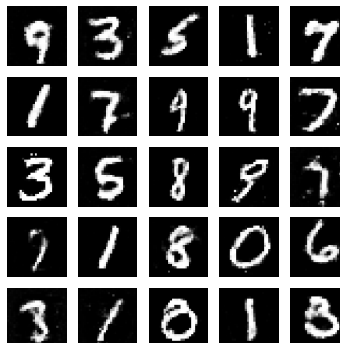
Generated samples obtained after 10^4 and 10^5 epochs of training using nets D , G with architectures α_D , α_G :

$$\alpha_D = ((28 \times 28, 240, 240, 1), (\text{LeakyReLU}, \text{LeakyReLU}, \text{Sigmoid})),$$

$$\alpha_G = ((100, 1200, 1200, 28 \times 28), (\text{ReLU}, \text{ReLU}, \text{Tanh})).$$



(a) 10^4 epochs



(b) 10^5 epochs

REFERENCES

- [1] Ian 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*. Curran Associates, Inc., 2014. doi:10.48550/arXiv.1406.2661.
- [2] Maxime Vandegar. Generative adversarial networks. https://github.com/MaximeVandegar/Papers-in-100-Lines-of-Code/tree/main/Generative_Adversarial_Networks, 2020. Accessed: 2024/04/03.