

Artificial Creativity

Federico Matteoni

Intelligent Systems for Pattern Recognition

A.Y. 2022/2023

Mat. 530257

Abstract—This work is a survey on some common machine learning systems aimed at image generation that surfaced in the last years. These systems have recently been used by the general public in viral fashion to generate artistic content, ranging from artistic repositions of submitted photos to completely new concepts produced by text-based prompts. Although exceptional, these systems also raised some concerns regarding copyright and intellectual property.

I. INTRODUCTION

In January 2021, OpenAI labs presented DALL-E [5], a text-to-image generative discrete variational autoencoder model to the public with a readily available demo online. In that demo, users could simply type a prompt and, after less than few minutes, a small selection of artistic interpretations of that prompt would be presented. While not being the first text-to-image machine learning model to be presented, it can be argued that DALL-E is the first that had a big impact on the public: in the following months many examples of images generated with that tool would be posted on social media, and many more sophisticated models would surface, taking advantage of the increasing interest in the image generation and *artificial creativity* field.

Many of the most recent and successful of these systems stem from the "Diffusion Probabilistic Models" [6] family of machine learning approaches, while some of the earliest examples are Generative Adversarial Networks [1] with alternative techniques for the sample generation [3] compared to the classical GAN approaches.

II. MODELS

Image generation, and broadly speaking every kind of data generated through machine learning, is achieved with a family of models called "Generative Models". These models, in contrast with the family of "Discriminative Models", aim to generate new samples of data after learning its distribution, while the goal of discriminative models is to distinguish between different samples of data. Mathematically speaking, given a set of samples X and the corresponding labels Y , a generative model learns $p(X, Y)$ while a discriminative model learns $p(X | Y)$. Generative models can be explicit, when they directly learn $p(X)$, or implicit, if they learn how to sample from $p(X)$ without directly learning it.

Text-to-image generation can be achieved with many different kinds of generative models, but we will focus on the two main family of models that recently achieved the best results, in literature as well as with the general public: GANs and DPMs.

A. Generative Adversarial Networks

A Generative Adversarial Network (GAN) [1] is composed of two main parts:

- A Generator \mathcal{G} that learns to generate plausible data with the goal of confusing the discriminator
- A Discriminator \mathcal{D} that learns to distinguish between the real data and the fake data, penalizing the generator when it produces non-plausible samples

The trick that GANs employ to be able to sample data from complex, high-dimensional distributions, is to sample from a very simple distribution (usually a Gaussian, $\mathcal{N}(0, 1)$) and to train a differentiable distribution (usually a Neural Network, \mathcal{G}) to transform the sampled random noise into data from the target distribution:

- Sample $z \sim \mathcal{N}(0, 1)$
- Generate $x = \mathcal{G}(z)$
- Discriminate $\mathcal{D}(x) = \mathcal{D}(\mathcal{G}(z))$

Given real data x and random noise $z \sim \mathcal{N}(0, 1)$, the objective function of a GAN is C defined in (1)

$$C = \min_{\mathcal{G}} \max_{\mathcal{D}} (\mathbb{E}_x[\log \mathcal{D}(x)] - \mathbb{E}_z[\log(1 - \mathcal{D}(\mathcal{G}(z)))]) \quad (1)$$

This represents a hard two-player game, where the optimal solution resides in a saddle point that corresponds to the concurrent min max.

Most of the works in literature that employ GANs for data generation focus on developing different strategies that adapt the general architecture of these models to the interested data domain. Among the first examples of image generations that reached the general public is thispersondoesnotexist.com, a website released by nVidia which employs a particular GAN architecture, called StyleGAN, to generate highly realistic human portraits. Released in December 2018 with StyleGAN [3] and updated in December 2019 with StyleGAN2 [4], this project may represent the first public impact with state-of-the-art automatic image generation via machine learning models. The original motivation behind StyleGAN was to unbox and better understand the image synthesis process and the properties of the latent spaces employed by the GANs, to both better control the generative process and to better compare different generators between each other. By taking inspiration from the style-transfer literature [2], which studies techniques of merging the content of an image with the artistic style of another, the StyleGAN proposed a redesigned generative process: the input latent code (random noise) z is mapped into an intermediate latent space \mathcal{W} which is used to control

the generator \mathcal{G} at each layer together with added Gaussian noise. This is in contrast to traditional generators, which feed the input directly to the generating network thus using the latent data just in the input layer. StyleGAN uses latent data to control the generative process at each layer, with \mathcal{W} fed via adaptive instance normalization (AdaIN, (2)): it normalizes each feature map x_i separately given the style $y = (y_s, y_b)$ which is synthesized from \mathcal{W} via learned transformations.

$$\text{AdaIN}(x_i, y) = y_{si} \frac{x_i - \mu(x_i)}{\sigma(x_i)} + y_{bi} \quad (2)$$

A main difference from classical style transfer [2] is that the style y of the target image is computed from \mathcal{W} , the latent data, instead of an example image.

The noise which is fed independently to each layer of the generative process is used to generate stochastic detail in the final images: these details range from the placement of the skin pores to other details like freckles and hair. By using random noise at each layer the network does not need to use activations from the previous layer, thus reusing data from the only source of input available, to generate stochastic effects: this reduces the occurrence of repetitive patterns, which can be particularly perceivable in the case of human faces. According to style transfer literature, the style of an image is reliably encoded by spatially invariant statistics like channel-wise mean or variance etc., and in StyleGAN the globally defined style affects global effects like lightning, pose and the like, while locally defined noise added independently to each pixel is only used to generate stochastic variation: if the generator tried to control a global effect, e.g. pose, with the noise, this would lead to spatially inconsistent decisions and then penalization by the discriminator. So the network learns to use global and local channels appropriately without explicit guidance, fully employing the power of the adversarial process.

Besides this person does not exist, StyleGANs have been used in other online demos like this artwork does not exist and this cat does not exist: an example of the latter is Fig. 1.



Fig. 1. Example of cat portrait generated by a StyleGAN.

B. Diffusion Probabilistic Models

One of the main limitation in machine learning, and in generative models in particular, is the tractability of the probability distributions involved. Learning, sampling and evaluation are often still analytically or computationally intractable when using highly flexible families of models, thus requiring a tradeoff between these two conflicting goals: flexibility and tractability. Diffusion Probabilistic Models (DPM) [6] aim to achieve both, by taking inspiration from statistical physics.

DPMs are built by two main processes:

- A Diffusion Process which converts any complex data distribution into a simple one (e.g. a Gaussian)
- A Reverse Process which defines the generative model

The diffusion process gradually converts an input data distribution $q(x^{(0)})$ into a tractable distribution $p(y)$ by repeated application of a Markov Diffusion Kernel (3) where β is the diffusion rate

$$T_p(y | y'; \beta) \quad (3)$$

$$p(y) = \int T_p(y | y'; \beta) p(y') dy' \quad (4)$$

$$q(x^{(t)} | x^{(t-1)}) = T_p(x^{(t)} | x^{(t-1)}; \beta) \quad (5)$$

The forward trajectory corresponds to performing T steps of diffusion starting from the data distribution and is defined in (6), which exhibits clear Markovian dynamics showing a prior and a first-order transition

$$q(x^{(0...T)}) = q(x^{(0)}) \prod_{t=1}^T q(x^{(t)} | x^{(t-1)}) \quad (6)$$

The forward diffusion kernel $q(x^{(t)} | x^{(t-1)})$ for example may correspond to a Gaussian Diffusion into a Gaussian distribution with identity covariance, defined in (7)

$$q(x^{(t)} | x^{(t-1)}) = \mathcal{N}(x^{(t-1)} \sqrt{1 - \beta_t}, \mathbf{I} \beta_t) \quad (7)$$

The generative distribution is trained to describe the same trajectory of the forward process but in reverse, thus yielding the reverse process described in (8)

$$p(x^{(0...T)}) = p(x^{(T)}) \prod_{t=1}^T p(x^{(t-1)} | x^{(t)}) \quad (8)$$

During learning, and in the Gaussian diffusion kernel scenario taken as example, only the mean and covariance of the kernel need to be estimated:

- $f_\mu(x^{(t)}, t)$ defines the mean of the reverse Markov transitions
- $f_\sigma(x^{(t)}, t)$, defines the covariance of the reverse Markov transitions

Both can be defined via neural networks or other function fitting technique and the cost of these functions defines the computational cost of the overall algorithm. Those are used in the reverse diffusion kernel defined in (9)

$$p(x^{(t-1)} | x^{(t)}) = \mathcal{N}(f_\mu(x^{(t)}, t), f_\sigma(x^{(t)}, t)) \quad (9)$$

REFERENCES

- [1] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial networks, 2014.
- [2] X. Huang and S. J. Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. *CoRR*, abs/1703.06868, 2017.
- [3] T. Karras, S. Laine, and T. Aila. A style-based generator architecture for generative adversarial networks. *CoRR*, abs/1812.04948, 2018.
- [4] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. *CoRR*, abs/1912.04958, 2019.
- [5] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever. Zero-shot text-to-image generation, 2021.
- [6] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. *CoRR*, abs/1503.03585, 2015.