



3rd ML-INFN Hackathon: Advanced Level

INTRODUCTION TO GENERATIVE MODELS

GAN, VAE and Flow-based models

Matteo Barbetti (INFN-Firenze, University of Firenze) Francesco Vaselli (INFN-Pisa, Scuola Normale Superiore)



OUTLINE

1. What is generative modeling?

• Overview of state-of-the-art algorithms

2. Generative Adversarial Networks

Definition, drawbacks and improvements

3. Variational Autoencoders

Definition, probabilistic interpretation and known problems

4. Normalizing Flows

A family of methods for constructing flexible learnable probability distributions

5. Physics use-cases

How generative models are used in our field



•

WHAT IS GENERATIVE MODELING?

Overview of state-of-the-art algorithms



M. Barbetti (INFN-Firenze)

A look at academic trends



Data drawn from Google Scholar using Pold87/academic-keyword-occurrence

The interest of Scientific Community towards **Generative Models** has grown in recent years thanks to <u>Deep Learning</u> <u>developments</u>.

Usage and further improvement of this or the other algorithm are driven by how much it looks promising to cope a **specific task**.



Some applications

Generative modeling algorithms have been mainly developed for applications in **Computer Vision** to <u>create</u>, <u>transform</u> or <u>improve</u> a broad set of images composed of human faces, animals, landscapes, cartoons, sketches, photos and much more.

Generative models can be mainly used for:

- Data augmentation
 - Creation of new and never-seen instances according to the reference dataset
- Super resolution
 - Enhancing the resolution of an input image keeping its quality as high as possible
- Inpainting
 - Reconstruction of missing pixels in the input images keeping realism and consistency
- Denoising
 - Removing noise from input instances minimizing the loss of information
- Translation
 - Creation or transformation of images from a domain to another (i.e. text-to-image translation)



M. Barbetti (INFN-Firenze)

Text-to-image with DALL·E 2

"Teddy bears mixing sparkling chemicals as mad scientists as a 1990s Saturday morning cartoon"



Images taken from the OpenAI blog post



M. Barbetti (INFN-Firenze)

What I cannot create, I do not understand.

– Richard Feynman

M. Barbetti (INFN-Firenze)

A more general formulation

Let X and Y be respectively the set of data instances and the corresponding set of labels. Considering a Machine Learning <u>classification task</u>, the models that we can set up to face the problem can be broadly divided into <u>two main categories</u>:

- Discriminative model model of the conditional probability of the target Y, given an observation x, symbolically P(Y|X=x);
- Generative model model of the conditional probability of the observable X, given a target y, symbolically P(X|Y=y).



$$\begin{array}{c|c} \mathcal{D}_{\text{train}} = \{x_1, \dots, x_N\} \\ \text{with } x_i \sim \mathcal{P}(x) \end{array}$$

M. Barbetti (INFN-Firenze)

Machine-Learnt generative model



Figure stolen from the OpenAI blog post

3rd ML-INFN Hackathon: Advanced Level

09



10

M. Barbetti (INFN-Firenze)





M. Barbetti (INFN-Firenze)







10

M. Barbetti (INFN-Firenze)

Ready to use implementations

The main Machine Learning frameworks (Keras/TensorFlow/PyTorch) offer several tutorial to implement *vanilla* **generative modeling algorithms**, such as <u>DCGAN</u> or <u>CVAE</u> for the creation of digit images from the MNIST dataset.

Since they are quite recent algorithms, <u>high-level</u> solutions for Normalizing Flows are not yet available, unless through some **custom libraries** such as <u>FlowTorch</u> or <u>PZFlow</u>.

Some work is ongoing to provide a package for ready-to-use implementations of generative models for Physics applications:

- <u>TFGenModels</u> is based on TensorFlow 2 (work in progress);
- <u>TorchGen</u> is based on PyTorch (still to release).

Torch<mark>G</mark>en



2 GENERATIVE ADVERSARIAL NETWORKS

Definition, drawbacks and improvements



M. Barbetti (INFN-Firenze)

Introduction to GANs

Generative Adversarial Networks (GAN) rely on the simultaneous training of two neural nets:

- the discriminator network (D) is trained by a <u>classification task</u> to separate the generator output from the reference dataset;
- the generator network (G) is trained by a <u>simulation task</u> to reproduce the reference dataset trying to fake the discriminator.

The discriminator goal is to **maximize** the separation between the generator output and the real data, while the generator, <u>driven by the discriminator errors</u>, aims to **minimize** the differences with the data.

This framework corresponds to a **minimax two-players game**

The generative model can be thought of as analogous to a team of counterfeiters, trying to produce fake currency and use it without detection, while the discriminative model is analogous to the police, trying to detect the counterfeit currency.

– Ian J. Goodfellow et al.







Schematic representation





M. Barbetti (INFN-Firenze)

Mathematical definition

The original **loss function** proposed by Ian J. Goodfellow *et al.* is

$$\mathcal{L}_{\text{GAN}}(D,G) = \mathbb{E}_{x \sim P_{\text{true}}}[\log D(x)] + \mathbb{E}_{z \sim P_{z}}[\log(1 - D(G(z)))]$$

where D outputs the <u>probability</u> that its input comes from the reference dataset. The generator goal is that the discriminator makes a mistake, resulting with the following **minimax game**:



$$\min_{G} \max_{D} \mathcal{L}_{\text{GAN}}(D, G)$$

Theoretically a <u>unique solution</u> exists, with G recovering the reference data distribution and D equal to ½ everywhere.

M. Barbetti (INFN-Firenze)

Unstable training dynamics

Traditional GAN systems suffer from many issues mainly due to an unstable training procedure:

- the generator may collapse producing only a single sample or a small family of very similar samples (problem named *mode collapse*);
- the two players may oscillate during training rather than converging to the equilibrium point;
- if **imbalance** between the two agents occurs, then the system is incapable of learning at all.



All these drawbacks result from the *vanishing gradient* problem, namely the <u>lack of information</u> for the update of the generator parameters. This is due to the <u>saturation</u> of the discriminator that is so good in distinguishing the origin of the two samples that no errors remain to the generator to <u>improve</u> the synthetic space.

M. Barbetti (INFN-Firenze)

Stabilizing the training procedure

Stabilizing the training procedure relies on preventing the <u>vanishing gradient</u> problem, that can be achieved redefining the loss function to implement the **Wasserstein distance** to avoid the discriminator <u>saturation</u>.

$$\mathcal{L}_{\text{WGAN}}(D,G) = \mathbb{E}_{x \sim P_{\text{true}}}[D(x)] - \mathbb{E}_{z \sim P_{z}}[D(G(z))]$$



with D parameterizing a **1-Lipschitz function**, that leads the minimax game to:

$$\min_{G} \max_{\|D\|_{L} \le 1} \mathcal{L}_{\mathrm{WGAN}}(D,G)$$

Several training strategies have been developed to ensure the discriminator lipschitzianity, that can be done **indirectly** (i.e. with WGAN-GP) or **directly** (i.e. with WGAN-ALP).

3rd ML-INFN Hackathon: Advanced Level

18

References for GANs

I. J. Goodfellow et al., "Generative Adversarial Nets", arXiv:1406.2661

- M. Mirza and S. Osindero, "Conditional Generative Adversarial Nets", <u>arXiv:1411.1784</u>
- A. Radford, L. Metz and S. Chintala, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", <u>arXiv:1511.06434</u>
- L. Metz, "Unrolled Generative Adversarial Networks", <u>arXiv:1611.02163</u>
- M. Arjovsky and L. Bottou, "Towards Principled Methods for Training Generative Adversarial Networks", arXiv:1701.04862
- M. Arjovsky, S. Chintala and L. Bottou, "Wasserstein GAN", <u>arXiv:1701.07875</u>
- I. Gulrajani et al., "Improved Training of Wasserstein GANs", <u>arXiv:1704.00028</u>
- M. G. Bellemare *et al.*, "The Cramer Distance as a Solution to Biased Wasserstein Gradients", <u>arXiv:1705.10743</u>
- D. Terjék, "Adversarial Lipschitz Regularization", <u>arXiv:1907.05681</u>



M. Barbetti (INFN-Firenze)

3. VARIATIONAL AUTOENCODERS

Definition, probabilistic interpretation and known problems

M. Barbetti (INFN-Firenze)



Introduction to VAEs

A *Variational AutoEncoder* (VAE) is an *autoencoder* (AE) whose training is regularized to avoid overfitting and ensure that the latent space has good properties that <u>enable generative process</u>. Also VAEs are based on the simultaneous training of two neural nets:

- the encoder network (e) defines a map from the reference space to the latent space with the goal of <u>build a low-dimensional representation</u> of the input;
- the decoder network (d) defines a map from the latent space to the synthetic space with the goal of maximize the matching between the output and the reference space.

Contrary to GANs, VAEs training doesn't rely on an adversarial procedure and the system goal is to minimize the distance between encoded-decoded data and the initial data (**reconstruction error**). Moreover, contrary to AEs, VAEs encode an input as a **distribution** (with good properties) over the latent space rather than a single point to enable generative capabilities.

Schematic representation

encoder training e process encoded vector (in latent space) decoder input d generation decoded content sampler process (reconstructed input / sampled vector (from latent space)

Scheme stolen from the Joseph Rocca's blog post

22

generated content)

Schematic representation





M. Barbetti (INFN-Firenze)

Mathematical details

The VAE training is driven by the following **loss function**:

 $\mathcal{L}_{\text{VAE}}(\theta, \phi) = \|x - d_{\theta}(z)\|^2 + \frac{D_{KL}\left(\mathcal{N}\left(\mu(x;\phi), \sigma(x;\phi)\right) \|\mathcal{N}(0,I)\right)}{2}$

where the first term is the *reconstruction error*, while the second one is a *regularization term* that forces the distribution induced by the encoder to have good properties for the **generative process**.





M. Barbetti (INFN-Firenze)

Probabilistic interpretation

Let x be observed variables, z latent variables and let p(x,z) be the <u>parametric model</u> of their joint distribution.

$$\mathcal{X} = \{x^{(i)}\}_{i=1}^{N} \quad \text{MLE} \quad \log p_{\theta}(\mathcal{X}) = \sum_{i=1}^{N} \log p_{\theta}(x^{(i)}) \quad \text{MLE} \quad \log p_{\theta}(\mathcal{X}) = \sum_{i=1}^{N} \log p_{\theta}(x^{(i)}) \quad \mathbb{E}$$

Typically, this marginal log-likelihood is <u>intractable</u> to compute. Exploiting the **Bayes' Theorem**, one should try to tackle the MLE maximizing its lower bound, named *Evidence Lower Bound* (ELBO):

$$\log p_{\theta}(x) = \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x,z) - \log q_{\phi}(z|x) \right] + D_{KL} \left(q_{\phi}(z|x) \| p_{\theta}(z|x) \right)$$
$$\geq \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x,z) - \log q_{\phi}(z|x) \right]$$
$$= \mathbb{E}_{q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right] - D_{KL} \left(q_{\phi}(z|x) \| p(z) \right) \equiv -\mathcal{L}_{\text{VAE}}$$

where q(z|x) is an *inference model* approximating the posterior p(z|x).

M. Barbetti (INFN-Firenze)

Known problems

Despite their strong probabilistic base, VAEs are and have been less used than GANs since the **limited capabilities** of the generative model provided. The major drawbacks of VAEs results from the <u>approximation of</u> <u>the maximum likelihood</u> previously mentioned:

- **blurry** or **fuzzy generated data** q(z|x) maps multiple x to the same encoding z, making p(x|z)<u>non-Gaussian</u> for which the MSE within the loss produces some "average" of p(x|z);
- *posterior collapse* (also called *optimization challenges* or *information preference*) the decoder is so powerful to minimize "alone" the reconstruction error <u>ignoring the latent space</u> and making collapsing the posterior q(z|x) to the prior p(z), namely $q(z|x) \approx p(z)$ and $D_{kl}(q||p) \approx 0$.

Several strategies have been developed to cope these drawbacks, most of which based on a <u>change</u> of the loss function (i.e. VLAE or beta-VAE), or on a <u>redesign</u> of the VAE architecture (i.e. VQ-VAE), or even on a <u>constraint</u> over the variational models usable as approximators (i.e. delta-VAE).



Input



Images stolen from the Enoch Kan's blog post

References for VAEs

- D. P. Kingma and M. Welling, "Auto-Encoding Variational Bayes", <u>arXiv:1312.6114</u>
- D. J. Rezende, S. Mohamed and D. Wierstra, "Stochastic Backpropagation and Approximate Inference in Deep Generative Models", <u>arXiv:1401.4082</u>
- X. Chen et al., "Variational Lossy Autoencoder", <u>arXiv:1611.02731</u>
- S. Zhao, J. Song and S. Ermon, "Towards Deeper Understanding of Variational Autoencoding Models", arXiv:1702.08658
- I. Higgins *et al.*, "beta-VAE: Learning Basic Visual Concepts with a Constrained Variational Framework", <u>ICLR 2017 proceeding</u>
- A. van den Oord, O. Vinyals and K. Kavukcuoglu, "Neural Discrete Representation Learning", arXiv:1711.00937
- J. Lucas *et al.*, "Understanding Posterior Collapse in Generative Latent Variable Models", <u>ICLR 2019</u> proceeding
- A. Razavi et al., "Preventing Posterior Collapse with delta-VAEs", <u>arXiv:1901.03416</u>

4. NORMALIZING FLOWS

A family of methods for constructing flexible learnable probability distributions



F. Vaselli (INFN-Pisa)

Why limit ourselves to just one sample?

GAN FLOW

Output: Image

Output: Distribution

from "Why I Stopped Using GAN — ECCV 2020"

The basic idea: change of variables

We define a transform *f* such that:

$$\begin{aligned} \mathbf{x} &= f(\mathbf{z}) \\ \mathbf{z} &= f^{-1}(\mathbf{x}) \end{aligned}$$

The two pdfs are related:

$$p_x(\mathbf{x})d\mathbf{x} = p_z(\mathbf{z})d\mathbf{z}$$

 $p_x(\mathbf{x}) = p_z(f^{-1}(\mathbf{x})) \det \left[\frac{d\mathbf{z}}{d\mathbf{x}}\right]$



The basic idea: image generation



 $p_{\mathbf{x}}(\mathbf{x})$

 $p_{\mathbf{z}}(\mathbf{z})$

The basic idea: complex transforms



from Lilian Weng

We need few building blocks

Task

Learn the **f(z)** to send $p_{\mathbf{z}}(\mathbf{z})$ into the (**unknown**) data distribution $p_{\mathbf{x}}(\mathbf{x})$

Pieces

- Basic distribution $p_{\mathbf{z}}(\mathbf{z})$, typically Gaussian
- Function called flow *f(z) invertible* and *differentiable*, with *tractable jacobian*

The usage is straightforward

Density evaluation

$$p_x(\mathbf{x}) = p_z(f^{-1}(\mathbf{x})) \det \left| \frac{d\mathbf{z}}{d\mathbf{x}} \right|$$

Sampling new data

- Sample from $p_{\mathbf{z}}(\mathbf{z})$ (Gaussian, trivial)
- Compute $\mathbf{x} = f(\mathbf{z})$ (fast)

The loss function is intuitive



 $\mathcal{L}(\phi) = -\mathbb{E}_{p_x^*(\mathbf{x})}[\log(p_z(f^{-1}(\mathbf{x}; \phi))) + \log(\det \mathbb{J}_{f^{-1}}(\mathbf{x}; \phi))]$

where ϕ are the parameters of f(z)

Splines are a smart choice for *f(z)*

Expressive

Admit analytical inverse, fast to invert AND evaluate

We use ML to learn the optimal disposition of points and derivatives

(Just one of the possible choices!)



Normalizing Flows are *powerful* GMs!

Efficient to sample from $p_{\mathbf{x}}(\mathbf{x})$

Efficient to evaluate $\, p_{\mathbf{x}}(\mathbf{x}) \,$

Highly expressive

Useful latent representation

Straightforward to train

Normalizing Flows are *flawed* GMs!

Computation of the Jacobian is hard

Not defined to work in a *discrete* contex!

Coupling layers address Jacobian complexity



Dequantization can be used to tackle discrete variables

Apply a gaussian smearing converting discrete data into a continuum



from arXiv:2001.11235

References for Normalizing Flows

- G. Papamakarios et al, "Normalizing Flows for Probabilistic Modeling and Inference", <u>arXiv:1912.02762</u>
 - An updated and comprehensive review of Normalizing Flows architectures and applications
- C. Durkan et al, "Neural Spline Flows", <u>arXiv:1906.04032</u>
 - Dear The implementation of Normalizing Flows through Rational Quadratic Splines
- Pyro Team, "Normalizing Flows-Introduction", <u>link</u>
 - A cool hands-on introduction to Normalizing Flow in the Pyro package (python). As of this date, there isn't a definitive package for NF neither for Tensorflow nor for Pytorch

5. PHYSICS USE-CASES

How generative models are used in our field

F. Vaselli (INFN-Pisa)



ML based Fast Simulation: GAN









from arXiv:1712.10321

ML based Fast Simulation: Flow



Anomaly Detection

Idea: Model pdf of signal and background using Normalizing Flows





 $p_{bkg} = p(x|M \pm \Delta)$

from LHC Olympics 2020



ML Generation for Mammograms





from arXiv:1807.03401

References for Physics use-cases

- M. Paganini *et al.*, "CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters with Generative Adversarial Networks", <u>arXiv:1712.10321</u>
- Claudius Krause, David Shih, "CaloFlow: Fast and Accurate Generation of Calorimeter Showers with Normalizing Flows", <u>arXiv:2106.05285</u>
- The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics, arXiv:2101.08320
- Stephen R. Green, Jonathan Gair, "Complete parameter inference for GW150914 using deep learning", <u>arXiv:2008.03312</u>
- D. Korkinof *et al.*, "High-Resolution Mammogram Synthesis using Progressive Generative Adversarial Networks", <u>arXiv:1807.03401</u>



F. Vaselli (INFN-Pisa)

CONCLUSIONS

Generative models are a powerful tool at our disposal, with widespread adoption in many Physics use-cases and convincing results

Different models have specific advantages and drawbacks

No readily available implementations for our problems, need to experiment



THANKS!

Any questions?

You can find us at:

- Matteo.Barbetti@fi.infn.it
- Francesco.Vaselli@pi.infn.it

49

BACKUP



Pedagogical explanation of GANs training

- 1. Minimax game near convergence: P_{gen} is similar to P_{true} and D is a partially accurate classifier;
- 2. D is trained to discriminate samples from data, <u>converging to optimality</u>;
- 3. The G update is driven by D that points to region that is more likely to be classified as data;
- **4.** After several steps of training, the networks will reach a point at which <u>both cannot improve</u> because D is unable to differentiate between the two distributions.



M. Barbetti (INFN-Firenze)

Conditional GANs

GANs provide an easy extension to the **conditional form**. Considering a training sample composed by <u>multi-variable</u> elements, we can imagine to split the variables into:

- variables whose distributions are the goal of the generator (Y)
- variables that simply conditionate the generator output (X)

The generator task remains that of reproducing <u>synthetic</u> <u>data</u> (Y), but now the conditional space (X) is joint to the latent space (R) to pursue this goal. On the other hand, the discriminator uses the **generator input/output** to infer the <u>data origin</u> (real or generated).



52

Kullback-Leibler divergence

The *Kullback–Leibler divergence* (also called *relative entropy*) is a measure of how one probability distribution is <u>different</u> from a second, reference probability distribution.

 $D_{KL}(P||Q) = \mathbb{E}_{x \sim P}[\log P(x)] - \mathbb{E}_{x \sim P}[\log Q(x)]$





M. Barbetti (INFN-Firenze)

Jensen-Shannon divergence

The *Jensen–Shannon divergence* is a method of measuring the <u>similarity</u> between two probability distributions.

$$D_{JS}(P||Q) = \frac{1}{2} D_{KL} \left(P||\frac{P+Q}{2} \right) + \frac{1}{2} D_{KL} \left(Q||\frac{P+Q}{2} \right)$$





M. Barbetti (INFN-Firenze)

Backpropagation in VAEs

In order to ensure that the error correctly *backpropagates* through the nets despite the presence of a <u>random sampling</u> in the halfway of the VAE architecture, a simple trick (called **reparameterization trick**) is used to make the <u>gradient descent</u> possible:

 $z = \mu(x; \phi) + \varepsilon \cdot \sigma(x; \phi)$ with $\varepsilon \sim \mathcal{N}(0, I)$



55