

Generative Adversarial Networks (GANs)

From Day1 why the instructor shouting for Transformation and Optimization

Why study generative modeling?

Generative models, and GANs in particular, enable machine learning to work with multi-modal outputs. For many tasks, a single input may correspond to many different correct answers, each of which is acceptable. Some traditional means of training machine learning models, such as minimizing the mean squared error between a desired output and the model's predicted output, are not able to train models that can produce multiple different correct answers.

From Day1 why the instructor shouting for transformation and Optimization

How do generative models work? How do GANs compare to others?

- Maximum likelihood estimation
 - To simplify the discussion somewhat, we will focus on generative models that work via the principle of maximum likelihood. Not every generative model uses maximum likelihood. Some generative models do not use maximum likelihood by default, but can be made to do so (GANs fall into this category).

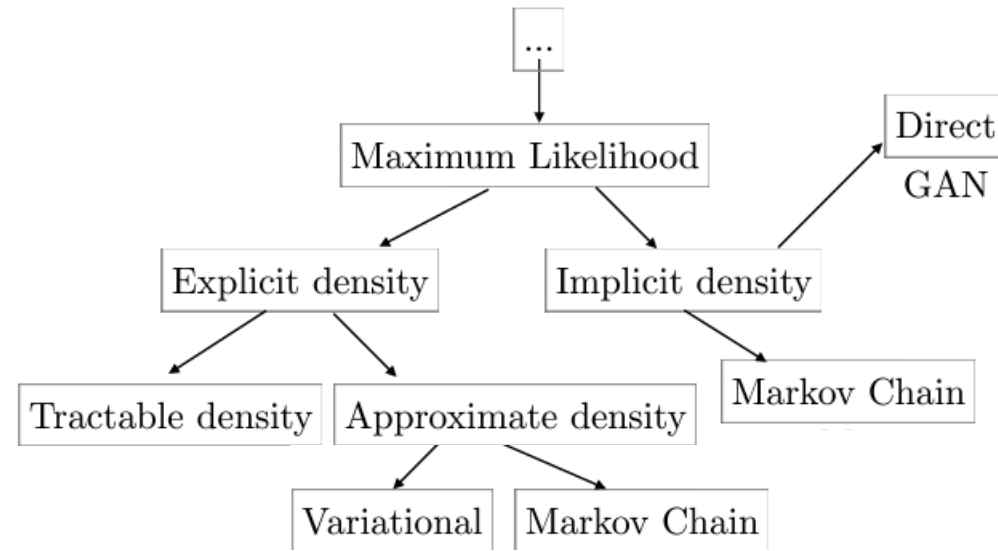
$$\begin{aligned}\theta^* &= \arg \max_{\theta} \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \\ &= \arg \max_{\theta} \log \prod_{i=1}^m p_{\text{model}}(\mathbf{x}^{(i)}; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{model}}(\mathbf{x}^{(i)}; \theta).\end{aligned}$$

From Day1 why the instructor shouting for transformation and Optimization

How do generative models work? How do GANs compare to others?

- Implicit density models

- Some models can be trained without even needing to explicitly define a density functions. These models instead offer a way to train the model while interacting only indirectly with p_{model} , usually by sampling from it.



From Day1 why the instructor shouting for transformation and Optimization

How do generative models work? How do GANs compare to others?

- Tractable explicit models
 - Fully visible belief networks are models that use the chain rule of probability to decompose a probability distribution over an n-dimensional vector \mathbf{x} into a product of one-dimensional probability distributions:

$$p_{\text{model}}(\mathbf{x}) = \prod_{i=1}^n p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$$

From Day1 why the instructor shouting for transformation and Optimization

Variational Autoencoders (VAEs)

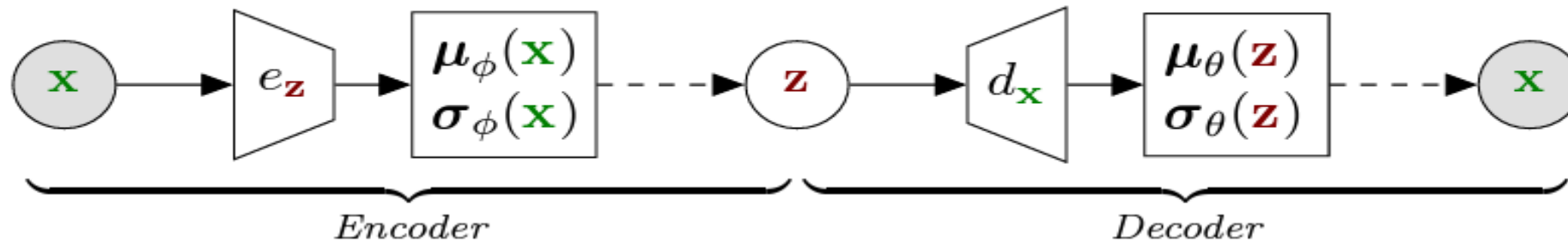
The mathematical basis of VAEs has relatively little to do with classical autoencoders (e.g. sparse autoencoders or denoising autoencoders). As mentioned earlier, our aim is to maximize the probability of each X in the training set under the entire generative process, according to:

$$P(X) = \int P(X|z;\theta)P(z)dz$$

From Day1 why the instructor shouting for transformation and Optimization

Variational Autoencoders (VAEs)

$\mathcal{N}(\cdot; \mu, \Sigma)$ denotes a multivariate Gaussian distribution with mean vector μ and covariance matrix Σ , $\text{diag}\{\cdot\}$ is the operator that forms a diagonal matrix from a vector by putting the vector entries on the diagonal, $\mathbf{0}_L$ is the zero-vector of size L , and \mathbf{I}_L is the identity matrix of size L . $p_{\theta_{\mathbf{x}}}(\mathbf{x})$ is a generic notation for a parametric pdf of the random variable \mathbf{x} , where $\theta_{\mathbf{x}}$ is the set of parameters. It is equivalent to $p(\mathbf{x}; \theta_{\mathbf{x}})$.



For simplicity of presentation and consistency across models, in the present review, $p_{\theta_{\mathbf{x}}}(\mathbf{x}|\mathbf{z})$ is assumed to be a Gaussian distribution with diagonal covariance matrix for all models; that is,

$$\begin{aligned} p_{\theta_{\mathbf{x}}}(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{x}; \mu_{\theta_{\mathbf{x}}}(\mathbf{z}), \text{diag}\{\sigma_{\theta_{\mathbf{x}}}^2(\mathbf{z})\}) \\ &= \prod_{f=1}^F p_{\theta_{\mathbf{x}}}(x_f|\mathbf{z}) = \prod_{f=1}^F \mathcal{N}(x_f; \mu_{\theta_{\mathbf{x}},f}(\mathbf{z}), \sigma_{\theta_{\mathbf{x}},f}^2(\mathbf{z})), \end{aligned}$$

From Day1 why the instructor shouting for transformation and Optimization

Denoising auto-encoder (DAE)

The Denoising Auto-Encoder (DAE) introduces a crucial modification to the traditional auto-encoder (AE) framework by incorporating a corruption process during training. This modification aims to enhance the robustness and generalization capabilities of the model by forcing it to learn from noisy or corrupted input data. The fundamental concept behind DAE lies in its ability to reconstruct the original clean input data \mathbf{x} from its corrupted counterpart $\hat{\mathbf{x}}$. The corruption process, represented by M_D , generates noisy versions of the input data, simulating real-world scenarios where the data may be corrupted by noise or other sources of interference.

$$\tilde{\mathbf{x}}^{(i)} \sim \mathcal{M}_D (\tilde{\mathbf{x}}^{(i)} | \mathbf{x}^{(i)})$$
$$L_{\text{DAE}}(\theta, \phi) = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - f_{\theta} (g_{\phi} (\tilde{\mathbf{x}}^{(i)})))^2$$

From Day1 why the instructor shouting for transformation and Optimization

Classical Auto-Encoder Models

Auto-encoders were originally proposed as artificial neural networks designed for dimension reduction. The fundamental concept behind classical auto-encoders involves constructing an architecture with input and output layers of the same dimension as the data while introducing intermediate layers of smaller dimensions to create a bottleneck. This bottleneck forces the network to learn a compressed representation of the input data, capturing its essential features. Classical auto-encoders have evolved over time to address various challenges and limitations inherent in the original design. One significant advancement of classical Auto-Encoders is the introduction of Denoising Auto-Encoders (DAE). DAEs are designed to reconstruct clean data from corrupted inputs, making them robust to noise and enhancing their generalization capabilities. Another notable development is the Stacked Convolutional Auto-Encoders (scAE). scAEs leverage convolutional layers to capture spatial dependencies and hierarchical features in the input data, making them particularly effective for imagerelated tasks. Additionally, Robust Deep Auto-Encoders (RDA), aim to enhance the robustness of auto-encoder models to outliers and anomalies in the input data.

From Day1 why the instructor shouting for transformation and Optimization

Diffusion models

Diffusion models have emerged as the state-of-the-art generation paradigm for a wide range of vision generation tasks. As recent challengers, autoregressive models share several characteristics with diffusion models. The methods are capable of generating diverse, high-quality samples, yet both also suffer from inefficient inference due to their iterative or sequential generation processes. Additionally, both approaches are likelihood-based and optimize objectives such as Negative Log-Likelihood (NLL) or Evidence Lower Bound (ELBO), making them relatively easy to train. Recent advances have focused on compressing visual content into latent spaces to improve the efficiency of high-resolution generation. Diffusion models employ a predefined forward process that gradually corrupts data and an iterative denoising process to recover samples. This iterative process retains spatial locality at each step, which benefits visual coherence. However, the necessity to corrupt the training data with Gaussian noise may potentially limit their scalability for understanding tasks. Additionally, while diffusion models operate primarily in continuous spaces, their discrete variants still lag behind in performance, posing challenges for multimodal applications like text generation. Autoregressive models, by contrast, inherently introduce unidirectional biases and discretization, which might not be ideal for visual tasks, partially explaining their lag behind diffusion models in recent benchmarks. However, autoregressive models offer flexibility in handling diverse modalities and combining generation with understanding. This adaptability aligns with the emerging trend of integrating autoregressive models with Large Language Models (LLMs) to create a unified framework for multimodal input-output tasks, bridging both generation and understanding capabilities.

From Day1 why the instructor shouting for transformation and Optimization

Masked Autoencoder.

Masked autoencoders (MAEs) are designed to learn robust data representations by randomly masking portions of the input data and training the model to reconstruct the missing content. Both MAEs and autoregressive models share similarities as they compress images into discrete sequences of visual elements and model these sequences. In fact, MAEs and autoregressive models inherit two leading paradigms from natural language processing (NLP): MAEs adopt the BERT-style encoder-decoder structure, while autoregressive models follow the GPT-Style decoder-only approach. Nevertheless, there are essential differences between them:

- 1) MAEs are trained by randomly masking visual tokens and reconstructing the missing part, while autoregressive models are trained to predict the next element in an ordered sequence;
- 2) MAEs utilize full attention, enabling each token to consider the entire surrounding contexts, whereas autoregressive models typically employ causal attention, restricting each token to focus only on previous ones;
- 3) During decoding, MAEs randomly generate multiple tokens in parallel, while autoregressive models generate tokens sequentially in a predefined order. These differences lead to distinct strengths: MAEs generally excel in representation learning and visual understanding tasks, and recent works like MaskGIT and MAGE also extend the Masked Image Modeling (MIM) approach to image generation.

From Day1 why the instructor shouting for transformation and Optimization