

Deep Generative Models

Ju Sun

Computer Science & Engineering
University of Minnesota, Twin Cities

April 30, 2026

Supervised learning

supervised learning: find **functional relationship** between input x and target y

- **Step 1** Gather training set $(x_1, y_1), \dots, (x_n, y_m)$
- **Step 2** Choose a NN with k neurons, so that there exist weights (w_1, \dots, w_k) to ensure $y_i \approx \{\text{NN}(w_1, \dots, w_k)\}(x_i), \forall i$
- **Step 3** Set up a loss function ℓ
- **Step 4** Find weights (w_1, \dots, w_k) to minimize the average loss

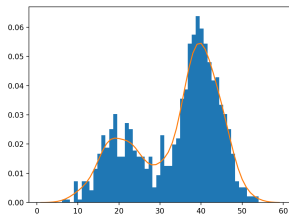
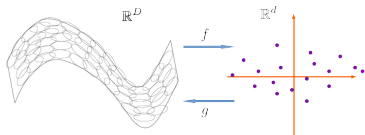
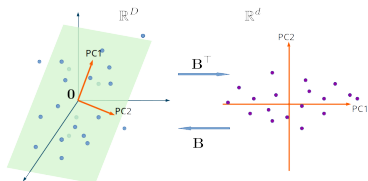
$$\frac{1}{m} \sum_{i=1}^m \ell[y_i, \{\text{NN}(w_1, \dots, w_k)\}(x_i)]$$

classification (y also called labels) or **regression**

object recognition, semantic segmentation, object detection, machine translation, image captioning, sentiment analysis, etc

Unsupervised learning

unsupervised learning: discover hidden **structures** in data, e.g., dimension reduction (subspace/manifold learning/sparse coding), clustering (e.g., k-means, spectral clustering), density estimation (e.g., GMM), etc

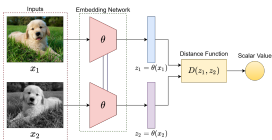


(Credit: <https://www.ecloudvalley.com/>)

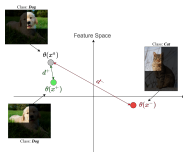
(Credit: <https://machinelearningmastery.com/>)

Self-supervised learning

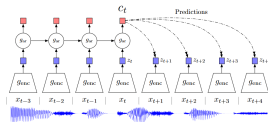
Marriage of unsupervised (no target) & supervised (contrived target) learning



joint embedding



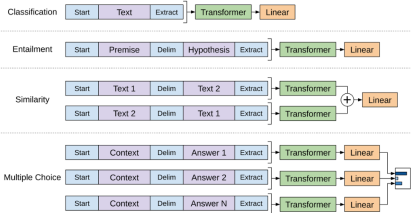
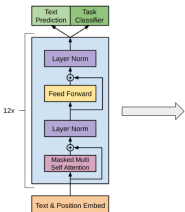
contrastive learning



sequential prediction

Pretraining: (transformer-based) decoder-only architectures pretrained on language model $\mathbb{P}[x^{(t+1)} | x^{(1)}, \dots, x^{(t)}]$

Finetuning: on task-specific supervised data

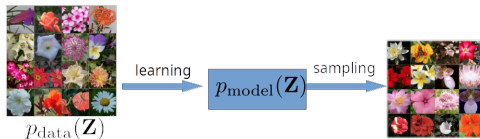


Generative learning/modeling: model distributions

notations: X, Y —RVs or events; \mathbf{x}, \mathbf{y} —samples of RVs; $\mathbb{P}[A]$ —probability of event A ; p : probability density function; θ in $p(\cdot; \theta)$ —parameters in parametrized density p

generative model: model $p(\mathbf{z})$, given iid samples $\{\mathbf{z}_i\}_{i=1}^n$

- generate new samples



- quantitative analysis: mean, variance, moments, probability (e.g., rare events)
- Bayesian learning/latent variable models: X —observations, Z —hidden variables

$$p(\mathbf{z} | \mathbf{x}) = p(\mathbf{x} | \mathbf{z})p(\mathbf{z})/p(\mathbf{x})$$

involving probability density functions

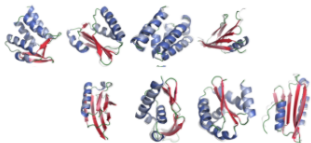
“Generative AI”



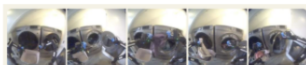
Text-2-Video
MovieGen, Meta



Text-2-Image
Stable Diffusion 3



Protein Generation
Huguet et al. 24



empty apartment dryer



batch fold shirts

Robot Action Model
Black et al. 24

Image credit: <https://neurips.cc/media/neurips-2024/Slides/99531.pdf>

Challenges

given samples $\{z_i\}_i$ drawn iid from $p_{\text{data}}(z)$

- **modeling:** which $p_{\text{model}}(z; \theta)$ gives good approximation?
- **computation:**
 - * parameter estimation: invoking *maximum likelihood estimation (MLE)*

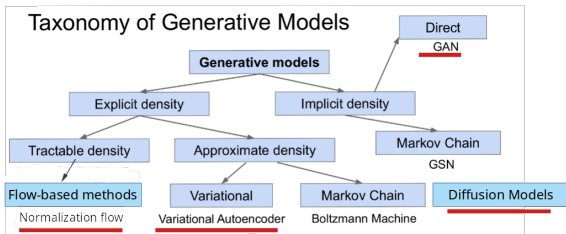
$$\max_{\theta} \prod_i p_{\text{model}}(z_i; \theta) \iff \max_{\theta} g(\theta) \doteq \sum_i \log p_{\text{model}}(z_i; \theta)$$

$g(\theta)$ can be highly nonconvex and hence hard to solve globally

- * quantitative analysis: many quantities, e.g., mean, variance, moments, likelihood, involve *high-dimensional integrals*
- **sampling:** *high-dimensional sampling* from even a **known** $p_{\text{model}}(z; \theta)$ can be expensive, e.g., Markov Chain Monte Carlo (MCMC)

classic methods: mixture models (e.g., GMM), kernel density estimation (Parzen window method)

Models we'll cover



(Credit: adapted from Stanford CS231N slides)

models focused on direct data generation

- Adversarial Generative Networks (GANs) [Goodfellow et al., 2014]
- Variational Autoencoder (VAE) [Kingma and Welling, 2013, Kingma and Welling, 2019]
- Diffusion Models & Flow-Matching Models [Sohl-Dickstein et al., 2015, Yang et al., 2022, Lipman et al., 2024]

models focused on density modeling & data generation

- Normalizing Flow [Rezende and Mohamed, 2015, Papamakarios et al., 2019]

Adversarial generative network (GAN)

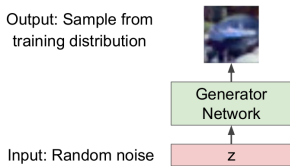
Variational autoencoder (VAE)

Normalizing flow

Diffusion models & Flow-Matching Models

Suggested reading

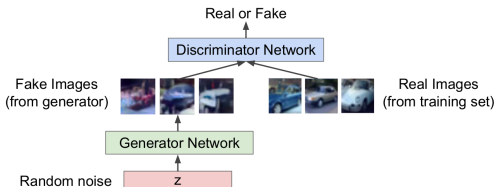
Learning via competition



(Credit: Stanford CS231N)

- **task:** given $\{x_i\}$, generate new samples from the distribution
- **idea:** map samples from a **simple** distribution to samples from the training distribution
- **how to:** measure the difference between **mapped/training** distributions

One solution: introduce a **biased** critic/discriminator



(Credit: Stanford CS231N)

- **Discriminator:** try to distinguish real (training) vs. fake (mapped) samples
- **Generator:** try to “fool” the discriminator by learning to generate realistic-looking images

Hope: the generator learns enough when the competition is stabilized

GAN objective

- **Discriminator**: try to distinguish real (training) vs. fake (mapped) samples
- **Generator**: try to “fool” the discriminator by learning to generate realistic-looking images

Hope: the generator learns enough when the competition is stabilized

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log \underbrace{D_{\theta_d}(\mathbf{x})}_{\text{discriminator output for real data}} + \mathbb{E}_{\mathbf{z} \sim p(Z)} \log[1 - \underbrace{D_{\theta_d}(G_{\theta_g}(\mathbf{z}))}_{\text{discriminator output for fake data}}]$$

- discriminator output: (0, 1) indicates likelihood of being real image (e.g., by passing sigmoid activations at output)
- discriminator wants to **maximize** the objective so that $D_{\theta_d}(\mathbf{x})$ is close to 1 and $D_{\theta_d}(G_{\theta_g}(\mathbf{z}))$ is close to 0
- generator wants to **minimize** the objective so that $D_{\theta_d}(G_{\theta_g}(\mathbf{z}))$ is close to 1

Training GANs: minimax optimization

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p(Z)} \log[1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))]$$

minimax (saddle point) optimization—way harder than minimization

heuristic algorithm: alternate between

- maximize wrt θ_d using **gradient ascent**

$$\max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p(Z)} \log[1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))]$$

- minimize wrt θ_g using **gradient descent**

$$\min_{\theta_g} \mathbb{E}_{\mathbf{z} \sim p(Z)} \log[1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))]$$

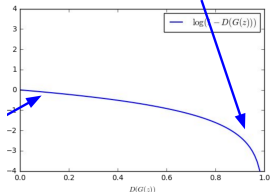
No guarantee this will work... [Razaviyayn et al., 2020]

Training GANs: the trick

minimize wrt θ_g using **gradient descent**

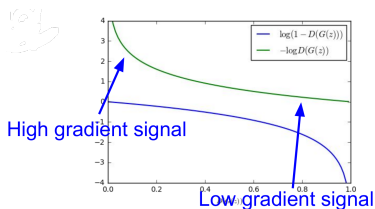
$$\min_{\theta_g} \mathbb{E}_{z \sim p(Z)} \log[1 - D_{\theta_d}(G_{\theta_g}(z))]$$

... but bit of numerical issue



(Credit: Stanford CS231N)

- grad. large when $D(G(z))$ is large, i.e., generator is already good
- grad. small when $D(G(z))$ is small, i.e., generator is not good, esp. at the initial stage—bad!



trick:

$$\min_{\theta_g} \mathbb{E}_{z \sim p(Z)} -\log D_{\theta_d}(G_{\theta_g}(z))$$

instead

GAN training pipeline

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D_{\theta_d}(x^{(i)}) + \log(1 - D_{\theta_d}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

end for

(Credit: Stanford CS231N)

- often $k > 1$ in practice to train the discriminator slightly faster
- overall still tricky and involves lots of tuning

Generate new samples

after training, use the generator to generate new samples

Output: Sample from training distribution



Generator Network

Input: Random noise

z

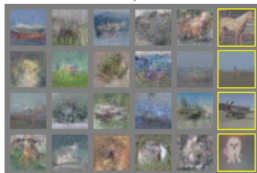
(Credit: Stanford CS231N)



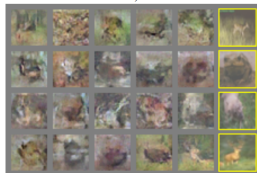
a)



b)



c)



d)

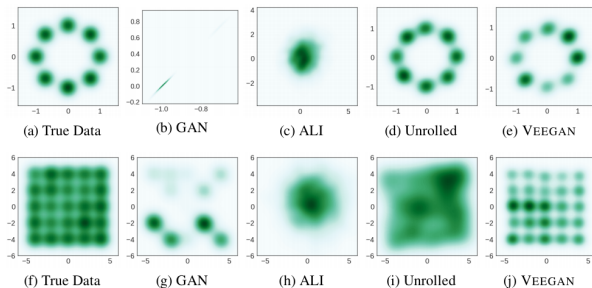
Figure 2: Visualization of samples from the model. Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set. Samples

(Credit: [Goodfellow et al., 2014])

What's wrong?

mode collapse: samples only generated from a subset of original support

Figure 2: Density plots of the true data and generator distributions from different GAN methods trained on mixtures of Gaussians arranged in a ring (top) or a grid (bottom).



(Credit: [Srivastava et al., 2017])

What really happens in the competition?

Lemma ([Goodfellow et al., 2014])

For any fixed G , the optimal discriminator D is $D_G^*(\mathbf{x}) = \frac{p_{\text{train}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x}) + p_g(\mathbf{x})}$.

Then

$$\begin{aligned} C(G) &= \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log D(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p(Z)} \log[1 - D(G(\mathbf{z}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log D_G^*(\mathbf{x}) + \mathbb{E}_{\mathbf{x} \sim p_g} \log[1 - D_G^*(\mathbf{x})] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log \frac{p_{\text{train}}(\mathbf{x})}{p_{\text{train}}(\mathbf{x}) + p_g(\mathbf{x})} + \mathbb{E}_{\mathbf{x} \sim p_g} \log \frac{p_g(\mathbf{x})}{p_{\text{train}}(\mathbf{x}) + p_g(\mathbf{x})} \end{aligned}$$

Theorem ([Goodfellow et al., 2014])

The global minimizer of $C(G)$ is at $p_g = p_{\text{train}}$.

Also, $C(G) = D_{\text{KL}}\left(P_{\text{train}} \parallel \frac{P_{\text{train}} + P_g}{2}\right) + D_{\text{KL}}\left(P_g \parallel \frac{P_{\text{train}} + P_g}{2}\right) - \log 4$, where D_{KL} denotes the KL-divergence $\text{KL}(P \parallel Q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$

Wasserstein GAN

Output: Sample from training distribution



Generator Network

Input: Random noise

z

(Credit: Stanford CS231N)

- **task:** given $\{x_i\}$, generate new samples from the distribution
- **idea:** map samples from a **simple** distribution to samples from the training distribution
- **how to:** measure the difference between **mapped/training** distributions

Jensen-Shannon (JS) divergence:

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

GAN tries to minimize $C(G) = 2D_{JS}(P_{\text{train}}\|P_g) + \text{const}$

Are there better measures for the difference?

Are all metrics created equal?

measure difference/distance between **distributions**

- **Total variation (TV) distance**

$$\text{TV}(P\|Q) = \sup_{A \in \mathcal{F}} |\mathbb{P}(A) - \mathbb{Q}(A)| \quad (\mathcal{F} : \text{sigma-algebra})$$

i.e., largest discrepancy of probabilities over all events A 's

- **Kullback-Leibler (KL) divergence** — asymmetric

$$\text{KL}(P\|Q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

- **Jensen-Shannon (JS) divergence** — symmetric

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

- **Earth-Mover (EM) or Wasserstein-1 distance (or Kantorovich–Rubinstein metric)**

$$W_1(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \|\mathbf{x} - \mathbf{y}\|$$

Earth mover distance/W-distance



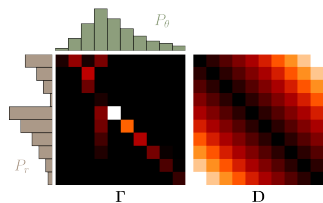
(Credit: <https://vincentherrmann.github.io/blog/wasserstein/>)

problem: move around stacks of earth P_r to match the shape/distribution of P_θ

goal: minimize the total movement effort, measured in **distance** \times **volume**

restrict P_θ and P_r to distributions

- let γ be a **transport plan**, Γ set of all plans.
- $\sum_x \gamma(x, y) = P_r(y) \forall y$: total move-in mass matches the target
- $\sum_y \gamma(x, y) = P_\theta(x) \forall x$: total move-out mass matches the original
- $\text{EMD}(P_r, P_\theta) = \inf_{\gamma \in \Gamma} \sum_{x, y} \|x - y\| \gamma(x, y) = \inf_{\gamma \in \Gamma} \mathbb{E}_{(x, y) \sim \gamma} \|x - y\|$



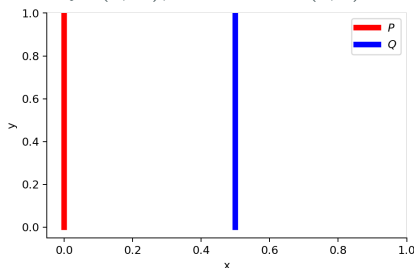
(Credit:

<https://vincentherrmann.github.io>)

Why Wasserstein distance?

$P : (0, W), W \sim \text{Uniform}(0, 1)$

$Q : (\theta, V), V \sim \text{Uniform}(0, 1)$

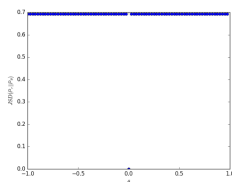
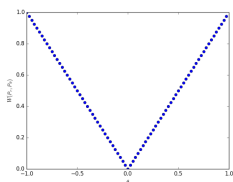


$$- \text{TV}(P||Q) = \begin{cases} 1 & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$- D_{\text{KL}} = \begin{cases} \infty & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$- D_{\text{JS}} = \begin{cases} \log 2 & \theta \neq 0 \\ 0 & \theta = 0 \end{cases}$$

$$- W_1(P, Q) = |\theta|$$



(Credit: [Arjovsky et al., 2017])

W-distance vs. JS-divergence

Wasserstein GAN (WGAN)

$$W_1(P, Q) = \inf_{\gamma \in \Gamma(P, Q)} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \gamma} \|\mathbf{x} - \mathbf{y}\|$$

hard to compute due to the minimization with the Γ constraint

Lemma (Kantorovich-Rubinstein duality)

$W_1(P, Q) = \sup_{f: \|f\|_L \leq 1} (\mathbb{E}_{\mathbf{x} \sim P} f(\mathbf{x}) - \mathbb{E}_{\mathbf{x} \sim Q} f(\mathbf{x}))$, where $\|f\|_L \leq 1$ denotes all functions with Lipschitz constant no greater than 1. (A function f is Lipschitz with a constant C if $\|f(\mathbf{x}) - f(\mathbf{y})\| \leq C \|\mathbf{x} - \mathbf{y}\|$ for all \mathbf{x}, \mathbf{y} .)

WGAN:

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} D_{\theta_d}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(Z)} D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \quad \text{s. t. } \|D_{\theta_d}\|_L \leq 1.$$

GAN:

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} \log D_{\theta_d}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(Z)} - \log[1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))]$$

Isn't it just a simple modification?

To train

$$\min_{\theta_g} \max_{\theta_d} \mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} D_{\theta_d}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(Z)} D_{\theta_d}(G_{\theta_g}(\mathbf{z})) \quad \text{s. t. } \|D_{\theta_d}\|_L \leq 1$$

the challenge is to enforce the Lipschitz constraint. Two heuristic:

- **weight clipping**: clip the discriminator θ_d into a predefined range $(-c, c)$.
But performance sensitive to choice of c
- **gradient penalty (GP)**: $\|D_{\theta_d}\|_L \leq 1 \implies \|\nabla_{\mathbf{x}} D\| \leq 1$, and max achieved when equality holds \implies encourage $\|\nabla_{\mathbf{x}} D\| = 1$

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} D_{\theta_d}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(Z)} D_{\theta_d}(G_{\theta_g}(\mathbf{z})) - \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{p}} (\|\nabla_{\tilde{\mathbf{x}}} D_{\theta_d}(\tilde{\mathbf{x}})\| - 1)^2 \right],$$

where $\tilde{\mathbf{x}} = t\mathbf{x} + (1-t)G_{\theta_g}(\mathbf{z})$ with $t \sim \text{uniform}(0, 1)$

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{\mathbf{x} \sim p_{\text{train}}} D_{\theta_d}(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(Z)} D_{\theta_d}(G_{\theta_g}(\mathbf{z})) - \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \tilde{p}} (\|\nabla_{\tilde{\mathbf{x}}} D_{\theta_d}(\tilde{\mathbf{x}})\| - 1)^2 \right],$$

where $\tilde{\mathbf{x}} = t\mathbf{x} + (1-t)G_{\theta_g}(\mathbf{z})$ with $t \sim \text{uniform}(0, 1)$

Algorithm 1 WGAN with gradient penalty. We use default values of $\lambda = 10$, $n_{\text{critic}} = 5$, $\alpha = 0.0001$, $\beta_1 = 0$, $\beta_2 = 0.9$.

Require: The gradient penalty coefficient λ , the number of critic iterations per generator iteration n_{critic} , the batch size m , Adam hyperparameters α, β_1, β_2 .

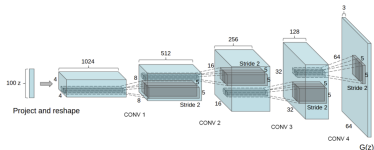
Require: initial critic parameters w_0 , initial generator parameters θ_0 .

```
1: while  $\theta$  has not converged do
2:   for  $t = 1, \dots, n_{\text{critic}}$  do
3:     for  $i = 1, \dots, m$  do
4:       Sample real data  $\mathbf{x} \sim \mathbb{P}_r$ , latent variable  $\mathbf{z} \sim p(\mathbf{z})$ , a random number  $\epsilon \sim U[0, 1]$ .
5:        $\tilde{\mathbf{x}} \leftarrow G_{\theta}(\mathbf{z})$ 
6:        $\hat{\mathbf{x}} \leftarrow \epsilon\mathbf{x} + (1 - \epsilon)\tilde{\mathbf{x}}$ 
7:        $L^{(i)} \leftarrow D_w(\tilde{\mathbf{x}}) - D_w(\mathbf{x}) + \lambda(\|\nabla_{\hat{\mathbf{x}}} D_w(\hat{\mathbf{x}})\|_2 - 1)^2$ 
8:     end for
9:      $w \leftarrow \text{Adam}(\nabla_w \frac{1}{m} \sum_{i=1}^m L^{(i)}, w, \alpha, \beta_1, \beta_2)$ 
10:   end for
11:   Sample a batch of latent variables  $\{\mathbf{z}^{(i)}\}_{i=1}^m \sim p(\mathbf{z})$ .
12:    $\theta \leftarrow \text{Adam}(\nabla_{\theta} \frac{1}{m} \sum_{i=1}^m -D_w(G_{\theta}(\mathbf{z})), \theta, \alpha, \beta_1, \beta_2)$ 
13: end while
```

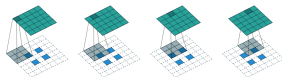
Deep convolutional GAN (DC-GAN)

both the generator G_{θ_g} and discriminator D_{θ_d} are deep convolutional networks [Radford et al., 2015]

the generator based on transposed (fractionally strided) convolutions



zero-padding, zero-dilution, then convolution



forward stride = 2

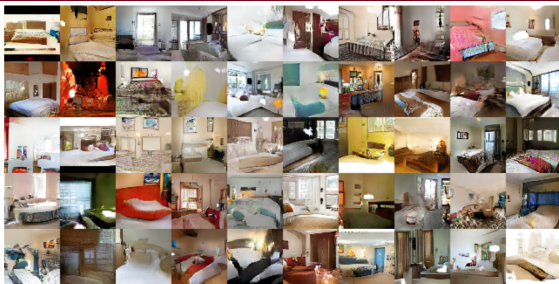
https://github.com/vdumoulin/conv_arithmetic

PyTorch implementation:

https://pytorch.org/tutorials/beginner/dcgan_faces_tutorial.html

DC-GAN results

image generation

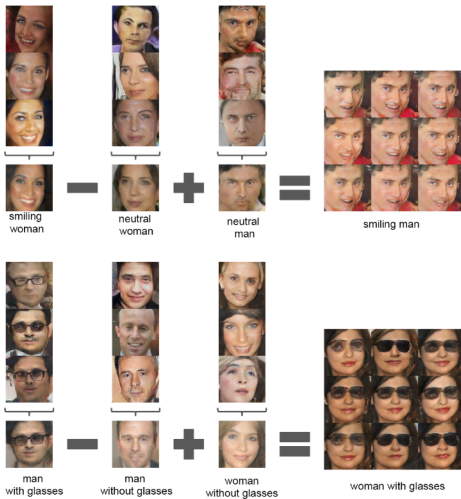


z interpolation



DC-GAN results

image arithmetic



Conditional GAN

augment additional info to both D and G inputs [Mirza and Osindero, 2014]

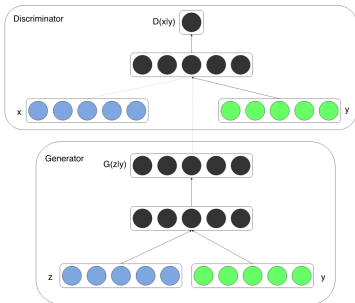






image + label (cond.)

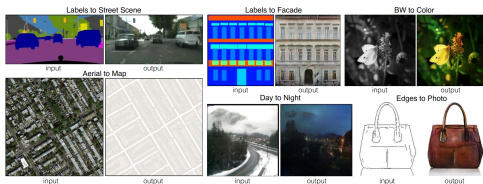


text generation + image (cond.)

	User tags + annotations	Generated tags
	montanha, trem, inverno, frio, people, male, plant life, tree, structures, transport, car	taxi, passenger, line, transportation, railway station, passengers, railways, signals, rail, rails
	food, raspberry, delicious, homemade	chicken, fattening, cookies, peanut, cream, cookie, house made, bread, biscuit, bakes
	water, river	creek, lake, along, near, river, rocky, treeline, valley, woods, waters
	people, portrait, female, baby, indoor	love, people, posing, girl, young, strangers, pretty, women, happy, life

Conditional GAN

image-to-image translation [Isola et al., 2016]



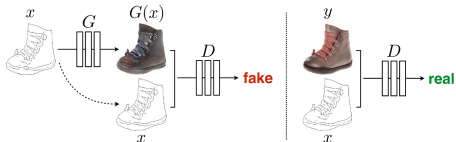
training data

Paired

x_i y_i



idea: conditional GAN + regression loss

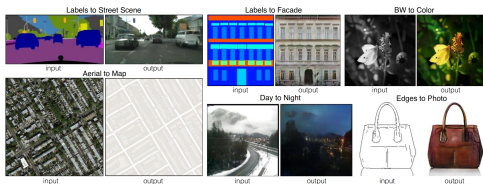


$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$

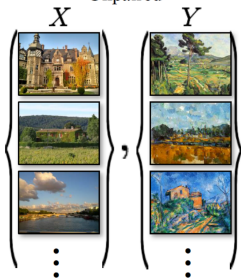
where

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,y} [\log D(x, y)] + \mathbb{E}_{x,z} [\log(1 - D(x, G(x, z)))]$$
$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z} [\|y - G(x, z)\|_1]$$

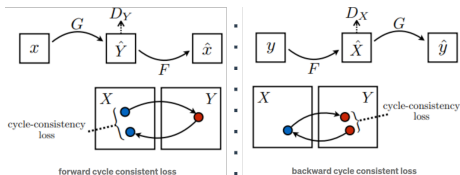
image-to-image translation **without** paired data [Zhu et al., 2017]



training data
Unpaired

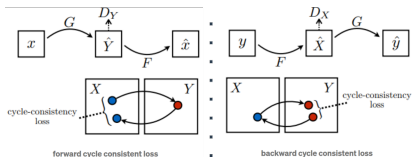


idea: match distributions with cycle consistency
 G : forward generator F : backward generator



$$\begin{aligned} x &\rightarrow G(x) \rightarrow F(G(x)) \approx x \\ y &\rightarrow F(y) \rightarrow G(F(y)) \approx y \end{aligned}$$

More on CycleGAN



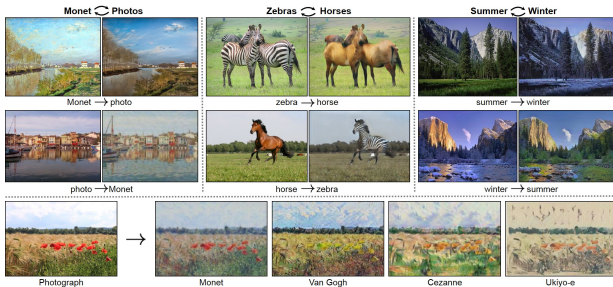
two generators \implies two GANs

$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cyc}}(G, F),$$

where $\mathcal{L}_{\text{cyc}}(G, F) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1].$

$$\begin{aligned} x &\rightarrow G(x) \rightarrow F(G(x)) \approx x \\ y &\rightarrow F(y) \rightarrow G(F(y)) \approx y \end{aligned}$$

\mathcal{L}_{cyc} looks familiar? autoencoder!



check out more <https://junyanz.github.io/CycleGAN/>

Progressive GAN

Toward **high-resolution generation**: progressively grow the generator output (and discriminator input) resolution as training goes on [Karras et al., 2017]

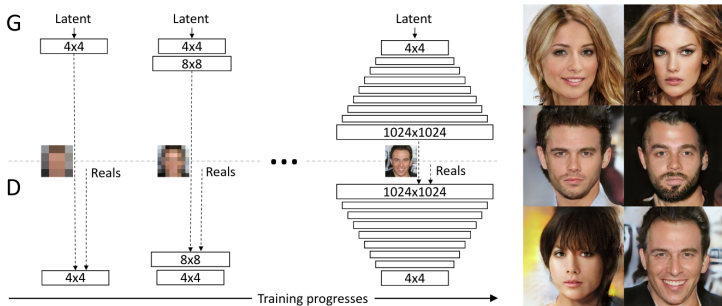


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $[N \times N]$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. On the right we show six example images generated using progressive growing at 1024×1024 .

Toward high-resolution generation [Brock et al., 2018]

LARGE SCALE GAN TRAINING FOR HIGH FIDELITY NATURAL IMAGE SYNTHESIS

Andrew Brock^{*†}
Heriot-Watt University
ajb5@hw.ac.uk

Jeff Donahue[†]
DeepMind
jeffdonahue@google.com

Karen Simonyan[†]
DeepMind
simonyan@google.com

ABSTRACT

Despite recent progress in generative image modeling, successfully generating high-resolution, diverse samples from complex datasets such as ImageNet remains an elusive goal. To this end, we train Generative Adversarial Networks at the largest scale yet attempted, and study the instabilities specific to such scale. We find that applying orthogonal regularization to the generator renders it amenable to a simple “truncation trick,” allowing fine control over the trade-off between sample fidelity and variety by reducing the variance of the Generator’s input. Our modifications lead to models which set the new state of the art in class-conditional image synthesis. When trained on ImageNet at 128×128 resolution, our models (BigGANs) achieve an Inception Score (IS) of 166.5 and Fréchet Inception Distance (FID) of 7.4, improving over the previous best IS of 52.52 and FID of 18.65.

1 INTRODUCTION



Figure 1: Class-conditional samples generated by our model.

Adversarial generative network (GAN)

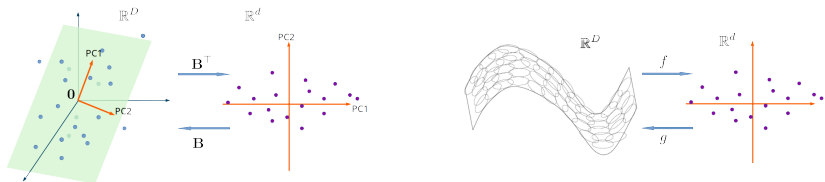
Variational autoencoder (VAE)

Normalizing flow

Diffusion models & Flow-Matching Models

Suggested reading

Quick summary of autoencoders and variants



	linear models	nonlinear models
autoencoder	$\min_B \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{B}^T \mathbf{x}_i)$ $\min_{B,A} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{A}^T \mathbf{x}_i)$	$\min_{V,W} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i))$
factorization	$\min_{B,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{z}_i)$	$\min_{V,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V(\mathbf{z}_i))$
sparse coding	$\min_{B,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, \mathbf{B}\mathbf{z}_i)$ $+\lambda \sum_{i=1}^m \Omega(\mathbf{z}_i)$	$\min_{V,Z} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V(\mathbf{z}_i))$ $+\lambda \sum_{i=1}^m \Omega(\mathbf{z}_i)$ $\min_{V,W} \sum_{i=1}^m \ell(\mathbf{x}_i, g_V \circ f_W(\mathbf{x}_i))$ $+\lambda \sum_{i=1}^m \Omega(f_W(\mathbf{x}_i))$

ℓ can be general loss functions other than $\|\cdot\|_2$

Ω promotes sparsity, e.g., $\Omega = \|\cdot\|_1$

Autoencoder vs. Variational autoencoder

Code in autoencoder

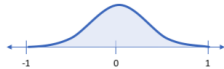
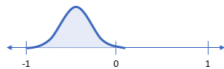
Code in variational autoencoder



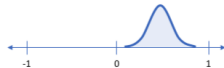
Smile (discrete value)



Smile (probability distribution)



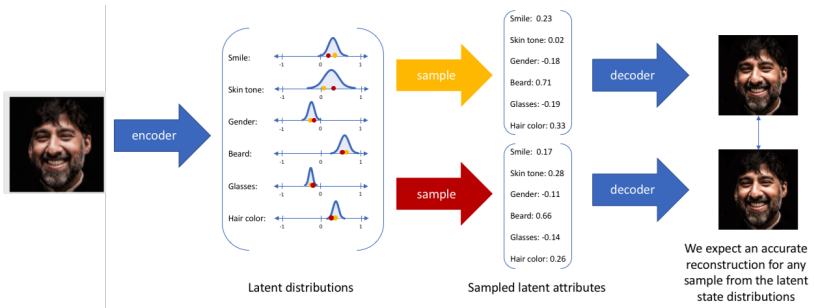
vs.



(Credit: <https://www.jeremyjordan.me/variational-autoencoders/>)

- Autoencoder maps each input to a **deterministic vector**
- Variational Autoencoder maps each input to a (parameterized) **distribution**

Variational autoencoder



(Credit: <https://www.jeremyjordan.me/variational-autoencoders/>)

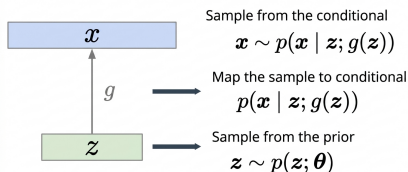
- **smoothness:** nearby codes tend to produce very similar reconstructions

Bayesian generation: to draw an \mathbf{x} from $p(\mathbf{x}, \mathbf{z})$, where $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$

- Step 1: draw a $\mathbf{z} \sim p(\mathbf{z})$
- Step 2: draw an $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z})$

Variational AE (VAE)

Model the generation via an **explicit density**



- $p(z; \theta)$, $p(x|z; \theta)$ are easy to sample from, e.g., $N(\mu, \Sigma)$
- the nonlinear mapping $g(z)$ allows expressive form of $p(x | z; g(z))$

GANs

Output: Sample from training distribution



Generator Network

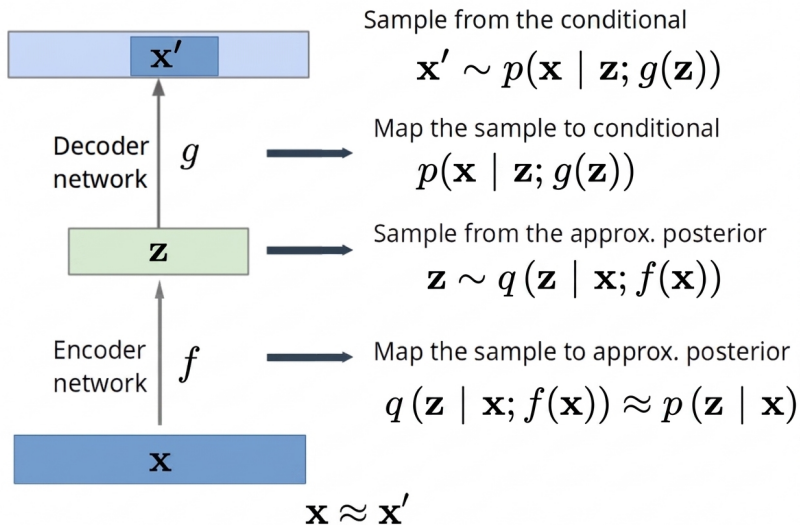
Input: Random noise

z

(Credit: Stanford CS231N)

- **task:** given $\{x_i\}$, generate new samples from the distribution
- **idea:** map samples from a **simple** distribution to samples from the training distribution
- **how to:** measure the difference between **mapped/training** distributions (JS, W-dist, etc)

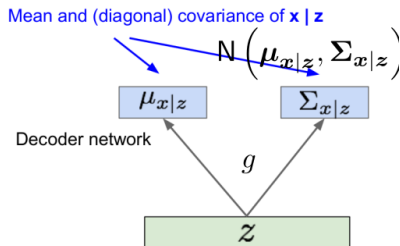
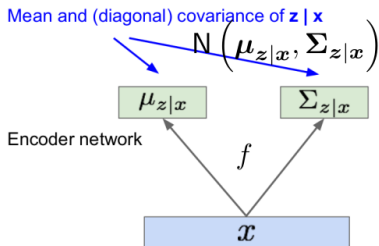
Augmenting the encoder



A VAE example

assume both $p(x | z)$ and $q(z | x)$ are multivariate Gaussian, i.e., $N(\mu, \Sigma)$

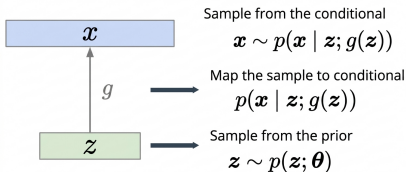
Since we're modeling probabilistic generation of data, encoder and decoder networks are probabilistic



(Credit: Stanford CS231N)

How to train it?

Model the generation via an **explicit density**



- $p(\mathbf{z}; \boldsymbol{\theta})$, $p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta})$ are easy to sample from, e.g., $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- the nonlinear mapping $g(\mathbf{z})$ allows expressive form of $p(\mathbf{x} | \mathbf{z}; g(\mathbf{z}))$

Assume m training samples $\mathbf{x}_1, \dots, \mathbf{x}_m$

$$\begin{aligned} \max_{\boldsymbol{\theta}, g} \prod_{i=1}^m p(\mathbf{x}_i; \boldsymbol{\theta}, g) &\iff \max_{\boldsymbol{\theta}, g} \sum_{i=1}^m \log p(\mathbf{x}_i; \boldsymbol{\theta}, g) \\ &\iff \max_{\boldsymbol{\theta}, g} \sum_{i=1}^m \log \int_{\mathbf{z}} p(\mathbf{x}_i | \mathbf{z}; g(\mathbf{z})) p(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \end{aligned}$$

(likely) intractable due to the integral

Monte Carlo sampling approx.? sample \mathbf{z}_j 's iid from $p(\mathbf{z}; \boldsymbol{\theta})$

$\int_{\mathbf{z}} p(\mathbf{x} | \mathbf{z}; g(\mathbf{z})) p(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \approx \frac{1}{K} \sum_{j=k}^K p(\mathbf{x} | \mathbf{z}_j; g(\mathbf{z}_j))$ —expensive and hard to converge

\mathbf{x} —observation, \mathbf{z} —latent variable

- $p(\mathbf{z})$ —prior
- $p(\mathbf{x} | \mathbf{z})$ —likelihood/conditional
- $p(\mathbf{z} | \mathbf{x})$ —posterior

want:

$$\max_{\theta, g} \sum_{i=1}^m \log p(\mathbf{x}_i; \theta, g)$$

Definition (Evidence lower bound (ELBO))

$\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}; q) \doteq \log p(\mathbf{x}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p(\mathbf{z} | \mathbf{x}))$ for any probability distribution q over \mathbf{z} (remember $D_{\text{KL}}(\cdot | \cdot) \geq 0$)

variational inference: find q so that the lower bound is tight as possible (i.e., $q(\mathbf{z} | \mathbf{x}) \approx p(\mathbf{z} | \mathbf{x})$ for all \mathbf{x}) but remains **tractable**

idea: restrict to a parameterized family $q(\mathbf{z} | \mathbf{x}; f(\mathbf{x}))$

(lots of other ideas in Bayesian inference, e.g., mean-field approximation. See, e.g., Chapter 19 of [Goodfellow et al., 2017])

How to train it?

need another identity

$$\begin{aligned}\log p(\mathbf{x}) - \text{KL}(q(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z} | \mathbf{x})) \\ = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x})} \log p(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))\end{aligned}$$

And maximize the right side instead

- maximizing $-D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) \| p(\mathbf{z}))$ ensures $q(\mathbf{z} | \mathbf{x})$ close to the prior $p(\mathbf{z})$
- maximizing $\mathbb{E}_{\mathbf{z} \sim q} \log p(\mathbf{x} | \mathbf{z})$ maximizes the likelihood of reproducing \mathbf{x} —minimizing reconstruction error
- overall, maximizing a lower bound to maximize the original

overall objective:

$$\max_{g, f} \sum_{i=1}^m \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z} | \mathbf{x}_i; f(\mathbf{x}_i))} \log p(\mathbf{x}_i | \mathbf{z}; g(\mathbf{x}_i)) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}_i; f(\mathbf{x}_i)) \| p(\mathbf{z}))$$

How to train it?

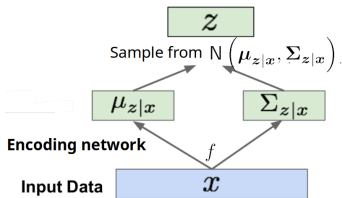
overall objective:

$$\max_{g,f} \sum_{i=1}^m \mathbb{E}_{z \sim q(z | \mathbf{x}_i; f(\mathbf{x}_i))} \log p(\mathbf{x}_i | z; g(\mathbf{x}_i)) - D_{\text{KL}}(q(z | \mathbf{x}_i; f(\mathbf{x}_i)) \| p(z))$$

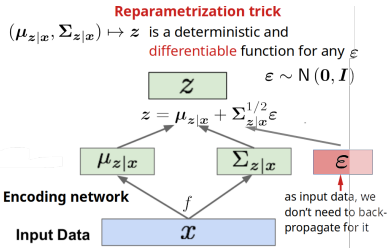
set $p(z) \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, $p(\mathbf{x}|z) \sim \mathcal{N}(\mathbf{u}_{\mathbf{x}|z}, \Sigma_{\mathbf{x}|z})$, $q(z|\mathbf{x}) \sim \mathcal{N}(\mathbf{u}_{z|\mathbf{x}}, \Sigma_{z|\mathbf{x}})$

build the computational graph for a single sample:

$(\mu_{z|\mathbf{x}}, \Sigma_{z|\mathbf{x}}) \mapsto z$ is a random function
and hence **not differentiable**

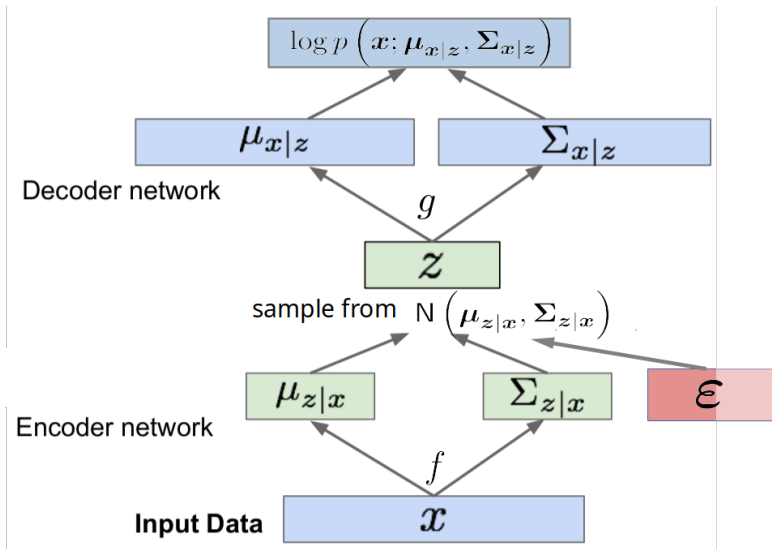


(Credit: adapted from Stanford CS231N)



(Credit: adapted from Stanford CS231N)

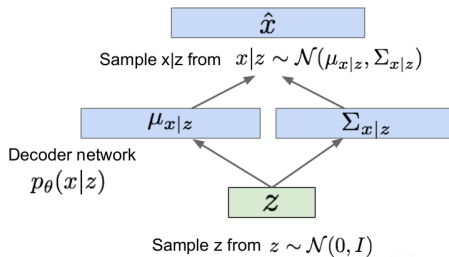
How to train it?



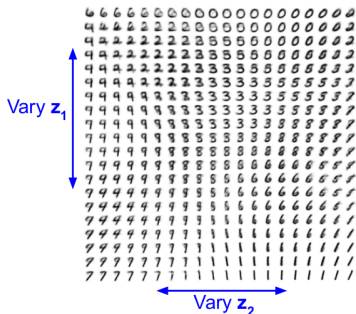
(Credit: adapted from Stanford CS231N)

Generate new samples

Use decoder network. Now sample z from prior!



Data manifold for 2-d z



(Credit: Stanford CS231N)

the coordinates of the codes potentially correspond to different physical properties (due to the diagonal covariance prior)

VAE objective:

$$\max_{g,f} \sum_{i=1}^m \mathbb{E}_{\mathbf{z} \sim q(Z | \mathbf{x}_i; f(\mathbf{x}_i))} \log p(\mathbf{x}_i | \mathbf{z}; g(\mathbf{x}_i)) - D_{\text{KL}}(q(Z | \mathbf{x}_i; f(\mathbf{x}_i)) \| p(Z))$$

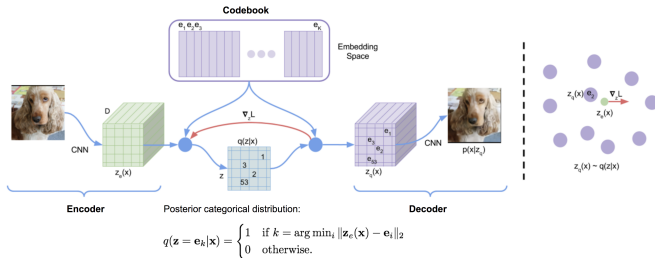
β -VAE objective [Higgins et al., 2017]:

$$\max_{g,f} \sum_{i=1}^m \mathbb{E}_{\mathbf{z} \sim q(Z | \mathbf{x}_i; f(\mathbf{x}_i))} \log p(\mathbf{x}_i | \mathbf{z}; g(\mathbf{x}_i)) - \beta D_{\text{KL}}(q(Z | \mathbf{x}_i; f(\mathbf{x}_i)) \| p(Z))$$

$\beta > 1$ to put more emphasis on the similarity of $q(Z | \mathbf{x}_i; f(\mathbf{x}_i))$ and $p(Z) \implies$ diagonal covariance of $p(Z)$ encourages decorrelation of coordinates in Z —disentangled representation



Vector quantization (VQ)-VAE



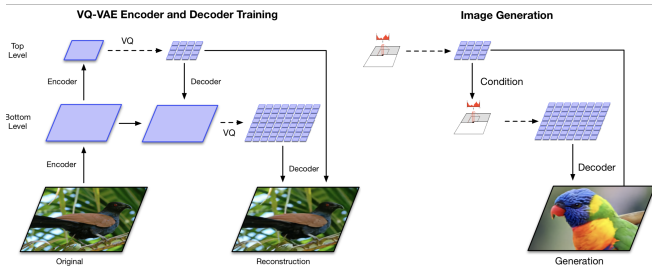
(Credit: [Van Den Oord et al., 2017])

- Finitely many latent codes, and hence discrete distribution
- Code assignment via nearest neighbor search
- Training via minimizing

$$L = \underbrace{\|x - D(e_k)\|_2^2}_{\text{reconstruction loss}} + \underbrace{\|sg[E(x)] - e_k\|_2^2}_{\text{VQ loss}} + \underbrace{\beta \|E(x) - sg[e_k]\|_2^2}_{\text{commitment loss}}$$

where sg is the stop_gradient operator

- Plus an auto-regressive prior



(Credit: [Razavi et al., 2019])

- Hierarchical VQ-VAE
- A prior learned on the discrete codebook
- In combination with self-attention enhanced autoregressive prior

Proof of the key equality

we'll omit the density probability for simplicity

$$\begin{aligned} & D_{\text{KL}}(q(Z|\mathbf{x})\|p(Z|\mathbf{x})) \\ &= \int q(\mathbf{z}|\mathbf{x}) \log \frac{q(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \int q(\mathbf{z}|\mathbf{x}) \log \frac{q(\mathbf{z}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{z}, \mathbf{x})} d\mathbf{z} \quad \text{Because } p(\mathbf{z}|\mathbf{x})=p(\mathbf{z}, \mathbf{x})/p(\mathbf{x}) \\ &= \int q(\mathbf{z}|\mathbf{x}) \left(\log p(\mathbf{x}) + \log \frac{q(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}, \mathbf{x})} \right) d\mathbf{z} \\ &= \log p(\mathbf{x}) + \int q(\mathbf{z}|\mathbf{x}) \log \frac{q(\mathbf{z}|\mathbf{x})}{p(\mathbf{z}, \mathbf{x})} d\mathbf{z} \quad \text{Because } \int q(\mathbf{z}|\mathbf{x}) d\mathbf{z}=1 \\ &= \log p(\mathbf{x}) + \int q(\mathbf{z}|\mathbf{x}) \log \frac{q(\mathbf{z}|\mathbf{x})}{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})} d\mathbf{z} \quad \text{Because } p(\mathbf{z}, \mathbf{x})=p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= \log p(\mathbf{x}) + D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p(\mathbf{x}|\mathbf{z}). \end{aligned}$$

So

$$\log p(\mathbf{x}) - D_{\text{KL}}(q(Z|\mathbf{x})\|p(Z|\mathbf{x})) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p(\mathbf{x}|\mathbf{z}) - D_{\text{KL}}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})).$$

Adversarial generative network (GAN)

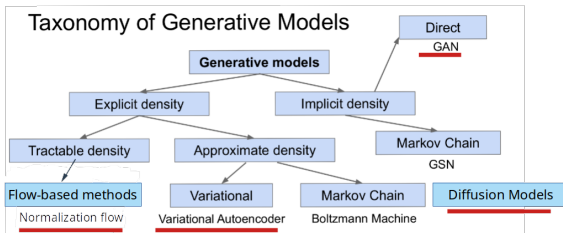
Variational autoencoder (VAE)

Normalizing flow

Diffusion models & Flow-Matching Models

Suggested reading

Modeling density directly



(Credit: adapted from Stanford CS231N slides)

Classical ideas:

- (Gaussian) Mixture models
- non-parametric methods, e.g., kernel density estimation

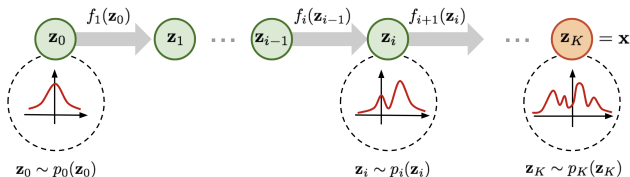
Change of variable theorem in probability

Consider a random variable Z and its density $\pi(\mathbf{z})$, what is the density of $\mathbf{x} = f(\mathbf{z})$, suppose that f is invertible?

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

where $\frac{df^{-1}(\mathbf{x})}{d\mathbf{x}}$ is the Jacobian of the inverse function $f^{-1}(\mathbf{x})$

Normalizing flow: idea



(Credit: <https://lilianweng.github.io/>)

$z_{i-1} \sim p_{i-1}(z_{i-1})$, $z_i = f_i(z_{i-1})$, thus $z_{i-1} = f_i^{-1}(z_i)$

$$p_i(z_i) = p_{i-1}(f_i^{-1}(z_i)) \left| \det \frac{df_i^{-1}}{dz_i} \right|$$

$$= p_{i-1}(z_{i-1}) \left| \det \left(\frac{df_i}{dz_{i-1}} \right)^{-1} \right|$$

According to the inverse func theorem.

$$= p_{i-1}(z_{i-1}) \left| \det \frac{df_i}{dz_{i-1}} \right|^{-1}$$

According to a property of Jacobians of invertible func.

$= \dots$

$$= p_0(z_0) \prod_{j=1}^i \left| \det \frac{df_j}{dz_{j-1}} \right|^{-1}$$

Normalizing flow: ideas

$$\begin{aligned} p_i(\mathbf{z}_i) &= p_0(\mathbf{z}_0) \prod_{j=1}^i \left| \det \frac{d\mathbf{f}_j}{d\mathbf{z}_{j-1}} \right|^{-1} \implies p(\mathbf{x}) = p_K(\mathbf{z}_K) = p_0(\mathbf{z}_0) \prod_{j=1}^K \left| \det \frac{d\mathbf{f}_j}{d\mathbf{z}_{j-1}} \right|^{-1} \\ &\implies \log p(\mathbf{x}) = \log p_0(\mathbf{z}_0) - \sum_{j=1}^K \log \left| \det \frac{d\mathbf{f}_j}{d\mathbf{z}_{j-1}} \right| \end{aligned}$$

So we can do maximum likelihood inference directly:

$$\max_{\boldsymbol{\theta}} \frac{1}{N} \sum_{\ell=1}^N \log p(\mathbf{x}_{\ell})$$

f 's are parametrized by DNNs with shared weights

Key challenges:

- ensure that f is invertible so that $\det \frac{d\mathbf{f}_j}{d\mathbf{z}_{j-1}}$ does not vanish
- computational tractability due to \det

Invertible maps in normalizing flows

Layer	Forward $f(x)$	Jacobian J	$\log \det J $
ActNorm	$z = s \odot x + b$	$\text{diag}(s)$	$\sum_{j=1}^D \log s_j $
Inv. 1×1 Conv	$z = Wx$	W	$\log \det W $
Affine Coupling	$\begin{bmatrix} \mathbf{x}_{1:d} \\ \mathbf{x}_{d+1:D} \odot e^{s(\mathbf{x}_{1:d}) + t(\mathbf{x}_{1:d})} \end{bmatrix}$	$\begin{bmatrix} I_d & \mathbf{0} \\ \frac{\partial \mathbf{z}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(e^s) \end{bmatrix}$	$\sum_j s(\mathbf{x}_{1:d})_j$

- For affine coupling, $s(\cdot)$ and $t(\cdot)$ are arbitrary neural networks.
- The lower-triangular structure of the Jacobian makes determinant calculation $O(D)$.

Adversarial generative network (GAN)

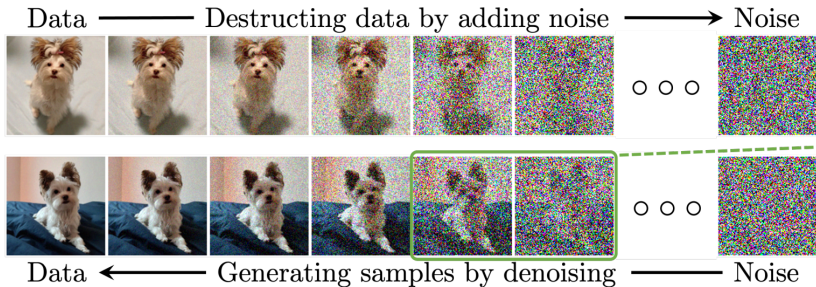
Variational autoencoder (VAE)

Normalizing flow

Diffusion models & Flow-Matching Models

Suggested reading

Image to noise, and noise back to image



(Credit: [Yang et al., 2022])

Forward diffusion: image \rightarrow noise

Reverse diffusion: noise \rightarrow image

Forward diffusion

Start from an image \mathbf{x}_0 and $\alpha_i \in (0, 1)$ for all i

- Step 1: $\mathbf{x}_1 = \sqrt{\alpha_1}\mathbf{x}_0 + \sqrt{1 - \alpha_1}\boldsymbol{\varepsilon}_1$ where $\boldsymbol{\varepsilon}_1 \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$
- Step 2: $\mathbf{x}_2 = \sqrt{\alpha_2}\mathbf{x}_1 + \sqrt{1 - \alpha_2}\boldsymbol{\varepsilon}_2$ where $\boldsymbol{\varepsilon}_2 \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$
- ...
- Step T : $\mathbf{x}_T = \sqrt{\alpha_T}\mathbf{x}_{T-1} + \sqrt{1 - \alpha_T}\boldsymbol{\varepsilon}_T$ where $\boldsymbol{\varepsilon}_T \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$

Now we have

$$\begin{aligned}\mathbf{x}_T &= \sqrt{\alpha_T}\mathbf{x}_{T-1} + \sqrt{1 - \alpha_T}\boldsymbol{\varepsilon}_T \\ &= \sqrt{\alpha_T} \left(\sqrt{\alpha_{T-1}}\mathbf{x}_{T-2} + \sqrt{1 - \alpha_{T-1}}\boldsymbol{\varepsilon}_{T-1} \right) + \sqrt{1 - \alpha_T}\boldsymbol{\varepsilon}_T \\ &= \sqrt{\alpha_T\alpha_{T-1}}\mathbf{x}_{T-2} + \sqrt{1 - \alpha_T\alpha_{T-1}}\boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \mathbf{I})\end{aligned}$$

Keep the induction, we obtain

$$\mathbf{x}_T = \sqrt{\prod_{t=1}^T \alpha_t} \mathbf{x}_0 + \sqrt{1 - \prod_{t=1}^T \alpha_t} \boldsymbol{\varepsilon} \text{ where } \boldsymbol{\varepsilon} \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$$

typically $\alpha_1 > \alpha_2 > \dots > \alpha_T$, and so $\mathbf{x}_T \sim \mathbf{N}(\mathbf{0}, \mathbf{I})$ as $T \rightarrow \infty$

Reverse diffusion

Assume the true prob. density in the forward diffusion $q(\mathbf{x}_t|\mathbf{x}_{t-1})$. If we also know $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we can reverse the process

Solution via simplification:

$q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ is approximated by

$$p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \doteq$$

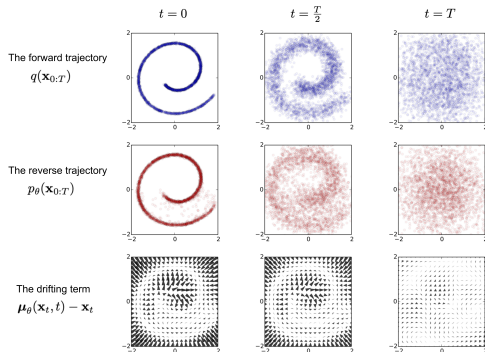
$N(\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$, where

$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$ and $\boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)$ are

learnable function parameterized

by DNNs

- Recall VAE?
- Training: minimize $\text{dist}(q(\mathbf{x}_0), p_{\theta}(\mathbf{x}_0))$



(Credit: [Sohl-Dickstein et al., 2015])

minimize $\text{dist}(q(\mathbf{x}_0), p_\theta(\mathbf{x}_0))$. If we can cross-entropy loss, then

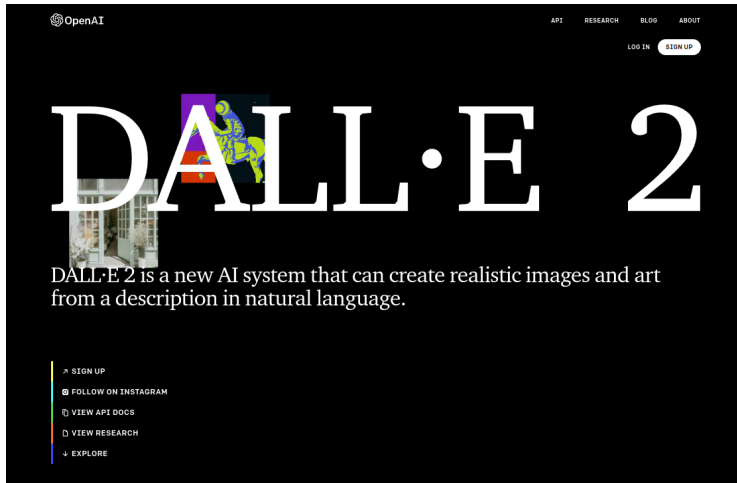
$$\begin{aligned} L_{\text{CE}} &= -\mathbb{E}_{q(\mathbf{x}_0)} \log p_\theta(\mathbf{x}_0) \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \right) \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\int q(\mathbf{x}_{1:T}|\mathbf{x}_0) \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} d\mathbf{x}_{1:T} \right) \\ &= -\mathbb{E}_{q(\mathbf{x}_0)} \log \left(\mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right) \\ &\leq -\mathbb{E}_{q(\mathbf{x}_{0:T})} \log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \\ &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] = L_{\text{VLB}} \end{aligned}$$

where VLB means evidence lower bound

After further rearrangement,

$$\begin{aligned} L_{\text{VLB}} &= \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}{p_\theta(\mathbf{x}_{0:T})} \right] \\ &= \mathbb{E}_q \left[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p_\theta(\mathbf{x}_T))}_{L_T: \text{constant}} + \right. \\ &\quad \left. \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))}_{L_{t-1}: \text{KL between Gaussians}} - \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \right] \end{aligned}$$

Diffusion models (and other models in action)



The image shows the OpenAI DALL·E 2 website landing page. At the top left is the OpenAI logo. The top right navigation bar includes links for API, RESEARCH, BLOG, and ABOUT, along with LOG IN and SIGN UP buttons. The main heading 'DALL·E 2' is displayed in large white serif font, with a colorful abstract image of a figure on a horse integrated into the letter 'A'. Below the heading, a smaller image of a modern interior is visible. The text below reads: 'DALL·E 2 is a new AI system that can create realistic images and art from a description in natural language.' At the bottom left, there is a vertical list of links: SIGN UP, FOLLOW ON INSTAGRAM, VIEW API DOCS, VIEW RESEARCH, and EXPLORE.

<https://openai.com/dall-e-2/>

Diffusion models (and other models in action)



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



an espresso machine that makes coffee from human souls, artstation



panda mad scientist mixing sparkling chemicals, artstation



a corgi's head depicted as an explosion of a nebula

<https://openai.com/dall-e-2/>

Diffusion models (and other models in action)

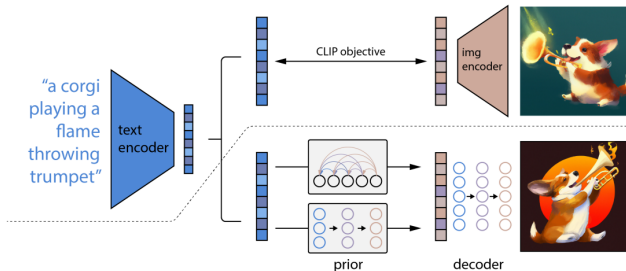
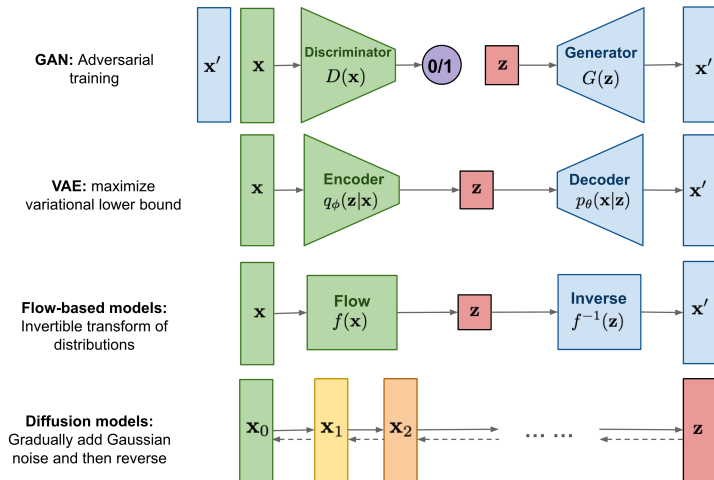


Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

<https://openai.com/dall-e-2/>

Summary of generative models



(Credit: <https://lilianweng.github.io/>)

Adversarial generative network (GAN)

Variational autoencoder (VAE)

Normalizing flow

Diffusion models & Flow-Matching Models

Suggested reading

Suggested reading

- CVPR 2018 tutorial on GANs [CVPR2018TutorialonGANs](#)
- NIPS 2016 Tutorial: Generative Adversarial Networks
<https://arxiv.org/abs/1701.00160>
- An Introduction to Variational Autoencoders [Kingma and Welling, 2019]
- Normalizing Flows for Probabilistic Modeling and Inference
<https://arxiv.org/abs/1912.02762>
- From GAN to WGAN <https://lilianweng.github.io/lil-log/2017/08/20/from-GAN-to-WGAN.html>
- From Autoencoder to Beta-VAE <https://lilianweng.github.io/lil-log/2018/08/12/from-autoencoder-to-beta-vae.html>
- Flow-based Deep Generative Models <https://lilianweng.github.io/lil-log/2018/10/13/flow-based-deep-generative-models.html#types-of-generative-models>
- What are Diffusion Models? <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

- [Arjovsky et al., 2017] Arjovsky, M., Chintala, S., and Bottou, L. (2017). **Wasserstein gan**. *arXiv:1701.07875*.
- [Brock et al., 2018] Brock, A., Donahue, J., and Simonyan, K. (2018). **Large scale gan training for high fidelity natural image synthesis**. *arXiv:1809.11096*.
- [Goodfellow et al., 2017] Goodfellow, I., Bengio, Y., and Courville, A. (2017). **Deep Learning**. The MIT Press.
- [Goodfellow et al., 2014] Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). **Generative adversarial networks**. *arXiv:1406.2661*.
- [Gulrajani et al., 2017] Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. (2017). **Improved training of wasserstein gans**. *arXiv:1704.00028*.
- [Higgins et al., 2017] Higgins, I., Matthey, L., Pal, A., Burgess, C., Glorot, X., Botvinick, M., Mohamed, S., and Lerchner, A. (2017). **beta-vae: Learning basic visual concepts with a constrained variational framework**. In *International Conference on Learning Representations, ICLR 2017*.

- [Isola et al., 2016] Isola, P., Zhu, J.-Y., Zhou, T., and Efros, A. A. (2016). **Image-to-image translation with conditional adversarial networks.** *arXiv:1611.07004*.
- [Karras et al., 2017] Karras, T., Aila, T., Laine, S., and Lehtinen, J. (2017). **Progressive growing of gans for improved quality, stability, and variation.** *arXiv:1710.10196*.
- [Kingma and Welling, 2013] Kingma, D. P. and Welling, M. (2013). **Auto-encoding variational bayes.** *arXiv:1312.6114*.
- [Kingma and Welling, 2019] Kingma, D. P. and Welling, M. (2019). **An introduction to variational autoencoders.** *arXiv:1906.02691*.
- [Lipman et al., 2024] Lipman, Y., Havasi, M., Holderrieth, P., Shaul, N., Le, M., Karrer, B., Chen, R. T., Lopez-Paz, D., Ben-Hamu, H., and Gat, I. (2024). **Flow matching guide and code.** *arXiv preprint arXiv:2412.06264*.
- [Mirza and Osindero, 2014] Mirza, M. and Osindero, S. (2014). **Conditional generative adversarial nets.** *arXiv:1411.1784*.

- [Papamakarios et al., 2019] Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2019). **Normalizing flows for probabilistic modeling and inference.** *Journal of Machine Learning Research*, 22(57):1-64, 2021.
- [Radford et al., 2015] Radford, A., Metz, L., and Chintala, S. (2015). **Unsupervised representation learning with deep convolutional generative adversarial networks.** *arXiv:1511.06434*.
- [Razavi et al., 2019] Razavi, A., van den Oord, A., and Vinyals, O. (2019). **Generating diverse high-fidelity images with vq-vae-2.** *arXiv:1906.00446*.
- [Razaviyayn et al., 2020] Razaviyayn, M., Huang, T., Lu, S., Nouiehed, M., Sanjabi, M., and Hong, M. (2020). **Non-convex min-max optimization: Applications, challenges, and recent theoretical advances.** *arXiv:2006.08141*.
- [Rezende and Mohamed, 2015] Rezende, D. J. and Mohamed, S. (2015). **Variational inference with normalizing flows.** *arXiv:1505.05770*.
- [Sohl-Dickstein et al., 2015] Sohl-Dickstein, J., Weiss, E. A., Maheswaranathan, N., and Ganguli, S. (2015). **Deep unsupervised learning using nonequilibrium thermodynamics.** *arXiv:1503.03585*.

- [Srivastava et al., 2017] Srivastava, A., Valkov, L., Russell, C., Gutmann, M. U., and Sutton, C. (2017). **Veegan: Reducing mode collapse in gans using implicit variational learning.** *arXiv:1705.07761*.
- [Van Den Oord et al., 2017] Van Den Oord, A., Vinyals, O., et al. (2017). **Neural discrete representation learning.** *Advances in neural information processing systems*, 30.
- [Yang et al., 2022] Yang, L., Zhang, Z., Song, Y., Hong, S., Xu, R., Zhao, Y., Shao, Y., Zhang, W., Cui, B., and Yang, M.-H. (2022). **Diffusion models: A comprehensive survey of methods and applications.** *arXiv:2209.00796*.
- [Zhu et al., 2017] Zhu, J.-Y., Park, T., Isola, P., and Efros, A. A. (2017). **Unpaired image-to-image translation using cycle-consistent adversarial networks.** *arXiv:1703.10593*.