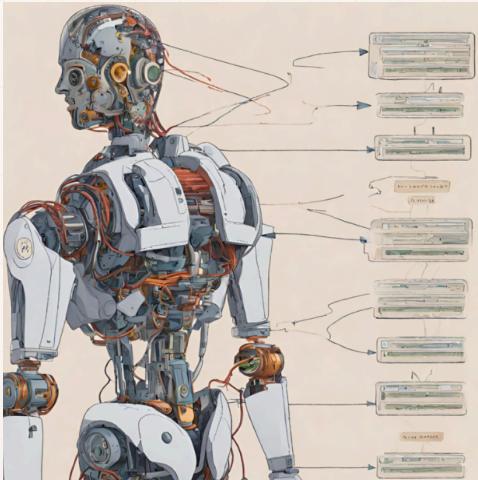


A TALE OF



VAEs, (VARIATIONAL AUTOENCODERS)
WGANs, (WASSERSTEIN GANS)
DDPMs (DENOISING DIFFUSION
PROBABILISTIC MODELS)

ASHWIN DE SILVA

Generative Modeling

Suppose that observed datum $x \sim p^*(x)$ "true distribution"

We wish to approximate $p^*(x)$ with a chosen model
 $p_\theta(x)$ with parameters $\theta \in \mathcal{H}$

Generative Modeling \equiv Process of searching for $\theta^* \in \mathcal{H}$ such that
 $p_{\theta^*}(x) \approx p^*(x)$

Generative Modeling

Dataset of observations $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\} \sim p^*(x)$ (iid)

Under the model $p_\theta(x)$,

$$\log p_\theta(\mathcal{D}) = \sum_{i=1}^N \log p_\theta(x^{(i)}) \equiv \text{log-likelihood of the data}$$

$$\max_{\theta} \sum_{i=1}^N \log p_\theta(x^{(i)})$$

Maximum Likelihood Estimate (MLE)

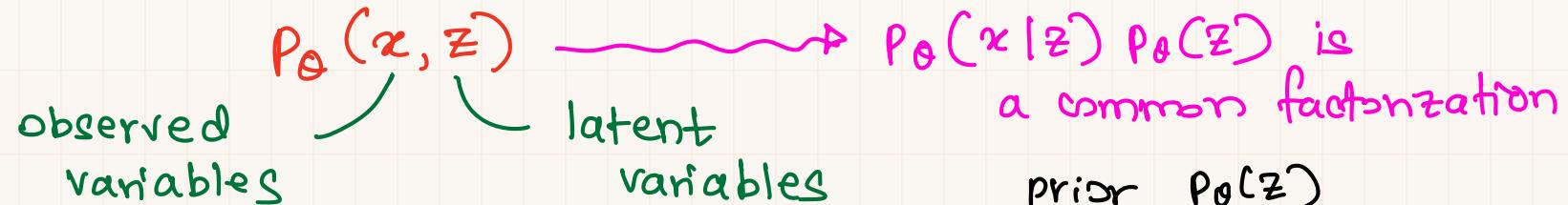
same as minimizing KL divergence between $p^*(x)$ and $p_\theta(x)$

can be solved by minibatch SGD

* In Bayesian inference, one may use MAP.

Latent Variable Models

Latent variables are the model variables that we don't observe.



$$P_\theta(x) = \int P_\theta(x, z) dz$$

marginal likelihood

↳ can lead to eg:-
highly expressive models

If z is discrete &
 $P_\theta(x|z)$ is Gaussian, then
 $P_\theta(x)$ is a GMM.

Latent Variable Models

Consider a LVM $P_{\theta^*}(x, z)$ with unknown θ^*

We have $x^{(1)}, \dots, x^{(N)} \sim P_{\theta^*}(x)$

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log P_{\theta}(x^{(i)})$$

* Cannot use SGD directly since taking $\nabla_{\theta} \log P_{\theta}(x)$

requires integrating over z ($\because P_{\theta}(x) = \int P_{\theta}(x, z) dz$)

We could use Bayes Theorem to obtain $P_{\theta}(x) = \frac{P_{\theta}(x, z)}{P_{\theta}(z|x)}$

But $P_{\theta}(z|x)$ is intractable to compute

Solution: approximate $P_{\theta}(z|x)$ by " $q_{\phi}(z|x)$ "

Evidence Lower Bound (ELBO)

$$\log P_\theta(x) \geq \mathbb{E}_{z \sim q(\cdot|x)} \left[\log \frac{P_\theta(x, z)}{q_\phi(z|x)} \right]$$

Instead of maximizing $\log P_\theta(x)$, maximize the ELBO:

$$\max_{\theta, \phi} \sum_{i=1}^N \mathbb{E}_{z \sim q(\cdot|x^{(i)})} \left[\log \frac{P_\theta(x^{(i)}, z)}{q_\phi(z|x^{(i)})} \right]$$

$$\equiv \min_{\theta, \phi} \sum_{i=1}^N \left\{ \text{KL}(q_\phi(z|x^{(i)}) \| P_\theta(z)) + \mathbb{E}_{z \sim q_\phi(\cdot|x^{(i)})} \left[\log P_\theta(x^{(i)}|z) \right] \right\}$$

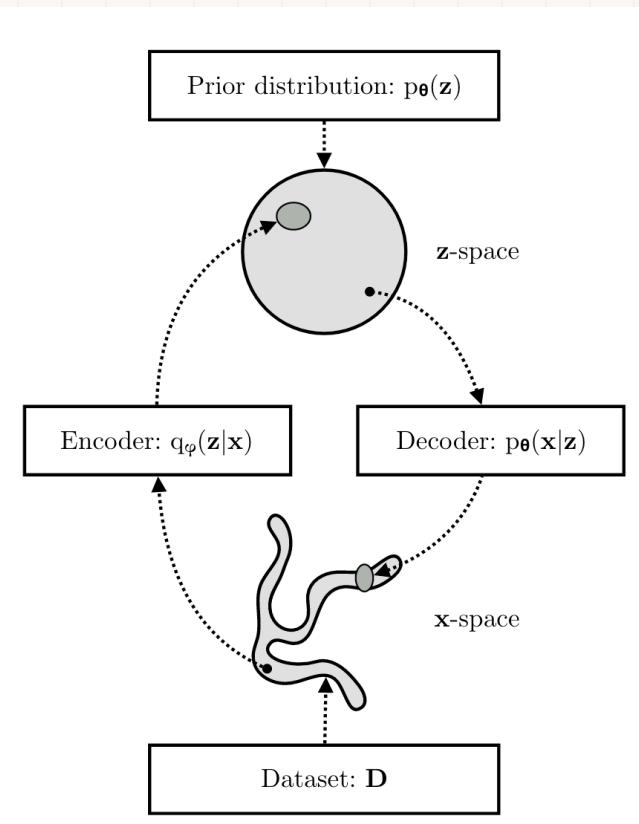
Variational Autoencoders (for Bernoulli data)

$$P(z) = \mathcal{N}(z; 0, I)$$

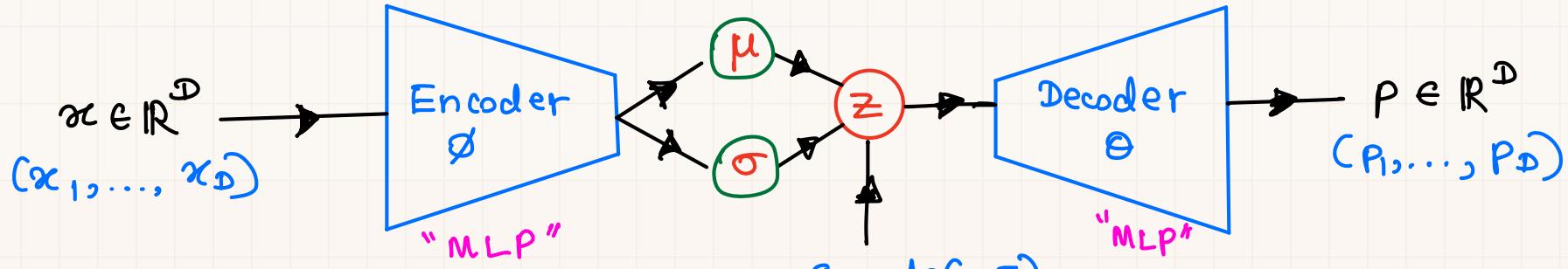
$$(p_1, \dots, p_D) = \text{Decoder } NN_{\theta}(z)$$

$$(\mu, \log \sigma) = \text{Encoder}_{\phi}(\alpha)$$

$$q_{\phi}(z|\alpha) = \mathcal{N}(z; \mu, \text{diag}(\sigma))$$



Variational Autoencoders (for Bernoulli data)



Loss function

$$\min_{\theta, \phi} \sum_{i=1}^N \left\{ \underbrace{\text{KL}\left(q_{\phi}(z|x^{(i)}) \parallel p_{\theta}(z)\right)}_{\text{KL between 2 gaussians}} - \mathbb{E}_{z \sim q_{\phi}(\cdot|x^{(i)})} \left[\underbrace{\log p_{\theta}(x^{(i)}|z)}_{\text{Binary Cross Entropy loss}} \right] \right\}$$

Wasserstein GAN

$$\text{Recall that } \max_{\theta} \sum_{i=1}^N \log p_{\theta}(x^{(i)}) \equiv \min_{\theta} \text{KL}(p_{\theta} \parallel p^*)$$

There're many measures of "distance" between distributions.

- KL-Divergence, J-S Divergence, ... (f-Divergences)
- Total Variation Distance • Hellinger Distance.
- Wasserstein Distance

$$W(p^*, p_{\theta}) = \inf_{\gamma \in \Pi(p^*, p_{\theta})} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|^2] \rightsquigarrow \begin{array}{l} \text{cost of the} \\ \text{optimal transport} \\ \text{plan} \end{array}$$

set of all
 joint distributions
 whose marginals are p^* & p_{θ}

Wasserstein GAN

Unlike other distances, $W(P^*, P_\theta)$

- is continuous everywhere
- is differentiable almost everywhere
- induces a weak topology in the space of distributions

} Desirable properties for a gradient based minimization.

$$W(P^*, P_\theta) = \inf_{\pi \in \Pi(P^*, P_\theta)} \mathbb{E}_{(x,y) \sim \pi} [\|x - y\|^2]$$

↑
↓
hard to evaluate

$$W(P^*, P_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P^*}[f(x)] - \mathbb{E}_{x \sim P_\theta}[f(x)]$$

↑ supremum over
1-Lipschitz functions

"Kantorovich-
Rubinstein
Duality"

Wasserstein GAN

- Suppose $z \sim \mathcal{N}(0, I)$. Then let g_θ be a generator
 - $x = g_\theta(z) \sim p_\theta(x)$
- Suppose $\{f_\phi\}$ be a family of functions parameterized by ϕ that are 1-Lipschitz

$$W(p^*, p_\theta) = \max_{\phi} \mathbb{E}_{x \sim p^*}[f_\phi(x)] - \mathbb{E}_z[f_\phi(g_\theta(z))]$$

Now we can minimize $W(p^*, p_\theta)$ w.r.t θ

$$\min_{\theta} W(p^*, p_\theta) = \min_{\theta} \max_{\phi} \mathbb{E}_{x \sim p^*}[f_\phi(x)] - \mathbb{E}_z[f_\phi(g_\theta(z))]$$

Wasserstein GAN

$$\min_{\Theta} \max_{\emptyset} \mathbb{E}_{x \sim p^*}[f_\phi(x)] - \mathbb{E}_z[f_\phi(g_\theta(z))]$$

"WGAN objective"

Algorithm 1 Training procedure of the Wasserstein GAN

Require: the learning rate α , the clipping parameter c , batch size m , number of iterations n_{critic} of the critic per generator iteration

Require: The initial critic parameters w_0 , the initial generator parameters θ_0

while θ has not converged **do**

for $t = 0, \dots, n_{\text{critic}}$ **do**

 Sample m real data $x^{(1)}, \dots, x^{(m)} \sim \mathbb{P}_r$

 Sample m latent realizations $z^{(1)}, \dots, z^{(m)} \sim p(z)$

$g_w \leftarrow \nabla_w \left[\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$

$w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$

$w \leftarrow \text{clip}(w, -c, c)$

 Sample m latent realizations $z^{(1)}, \dots, z^{(m)} \sim p(z)$

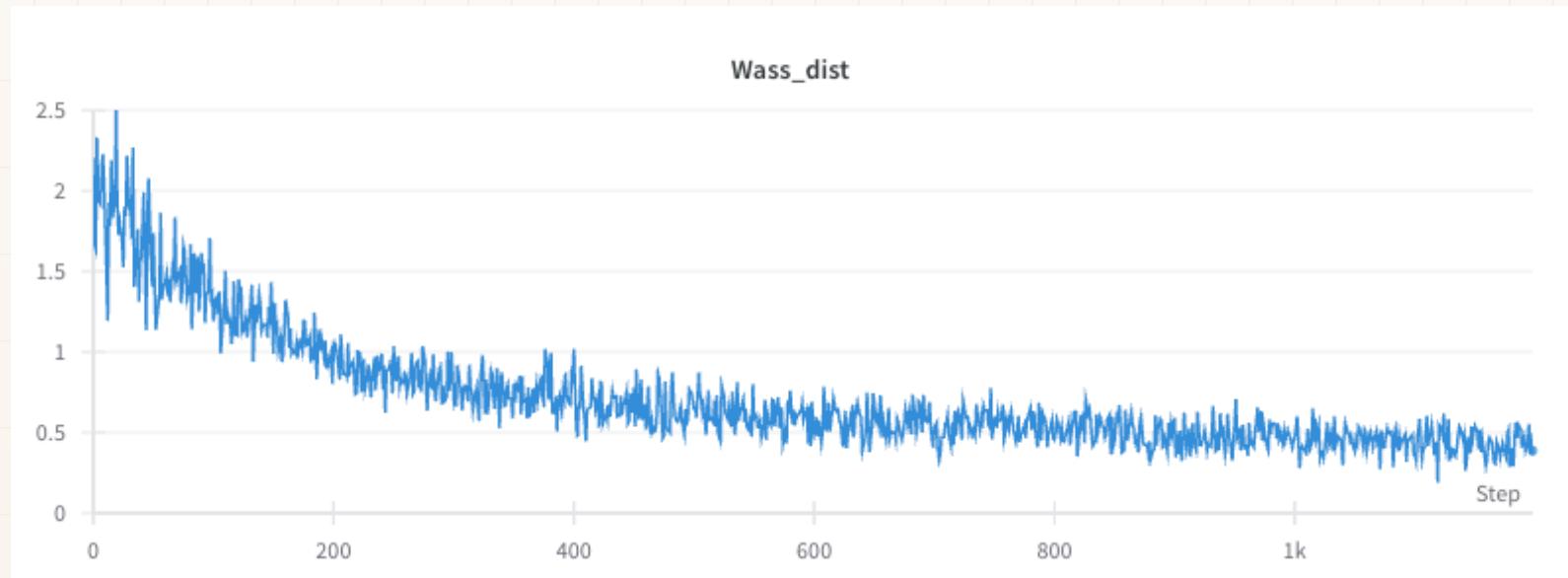
$g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$

$\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$

} "critic update"
S
} "Discriminator"
} "generator update"

Wasserstein GAN on MNIST

$W(p^*, p_\theta)$ estimate



* iterations

Denoising Diffusion Probabilistic Models

- Also a latent variable model $p_\theta(x_0, \underbrace{x_1, \dots, x_T}_{\text{latent variables}})$
observed \rightarrow
- $x_0 \sim q(x_0)$ true distribution
- let the joint $p_\theta(x_0, \dots, x_T)$ be a markov-chain (Reverse Process)

$$p_\theta(x_0, \dots, x_T) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t) \xrightarrow{\text{Normal}} \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

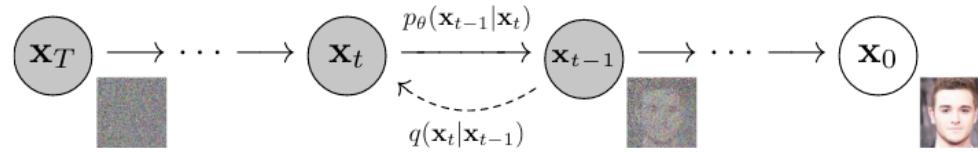
- Let the approximate posterior $q(x_1, \dots, x_T | x_0)$ also be a markov chain (forward process)

$$q(x_1, \dots, x_T | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

$\xrightarrow{\text{Normal}} \mathcal{N}(x_t; \sqrt{1 - \beta_t} x_t, \beta_t I)$

$\beta_t \sim \text{Noise Schedule}$

Denoising Diffusion Probabilistic Models



ELBO \Rightarrow

$$\mathbb{E}_{\theta}[-\log p_{\theta}(x_0)] \leq \mathbb{E}_q\left[-\log \frac{p_{\theta}(x_0, x_1, \dots, x_T)}{q(x_0, \dots, x_T | x_0)}\right] := L$$

"simplify"

$$h(\theta) = \mathbb{E}_{t, x_0, \varepsilon} \left[\| \varepsilon - \varepsilon_{\theta}(x_t, t) \|^2 \right]$$

Denoising Diffusion Probabilistic Models

$$h(\theta) = \mathbb{E}_{t, x_0, \epsilon} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|^2 \right]$$

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon$$

- ① Draw an image x_0
- ② Draw a noise $\epsilon \sim \mathcal{N}(0, I)$
- ③ Create a noisy image $x_t = x_0 + \epsilon$
- ④ Force network ϵ_θ to denoise x_t i.e. force ϵ_θ to predict ϵ .

Algorithm 1 Training

```
1: repeat
2:    $x_0 \sim q(x_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\epsilon \sim \mathcal{N}(0, I)$ 
5:   Take gradient descent step on
         $\nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\alpha_t} x_0 + \sqrt{1 - \alpha_t} \epsilon, t)\|^2$ 
6: until converged
```

But why?

Langevin Monte-Carlo, Score Matching & Denoising

(A high-level analogy)

Score function of probability density $p(x)$ is $s(x) = \nabla \log p(x)$

Langevin MC $\Rightarrow x_{t+1} = x_t - \eta \nabla \log p(x_t) + \sqrt{2\eta} \xi_t ; \xi_t \sim \mathcal{N}(0, I)$

Stationary Distribution of this SDE is $p(x)$

How to estimate $\nabla \log p(x)$?

$$L(\theta) = \mathbb{E}_{x \sim p} [\|\nabla \log p(x) - \varepsilon_\theta(x)\|^2]$$

{ approximately
equivalent to }

$$\text{"Denoising Loss"} \rightsquigarrow \mathbb{E}_{x \sim p, \varepsilon \sim \mathcal{N}(0, I)} [\|\varepsilon - \varepsilon_\theta(x + \varepsilon)\|^2]$$

\therefore By forcing a denoising objective, we let $\varepsilon_\theta(\cdot)$ approximate the score function $\nabla \log p(x)$

Sampling

Algorithm 2 Sampling

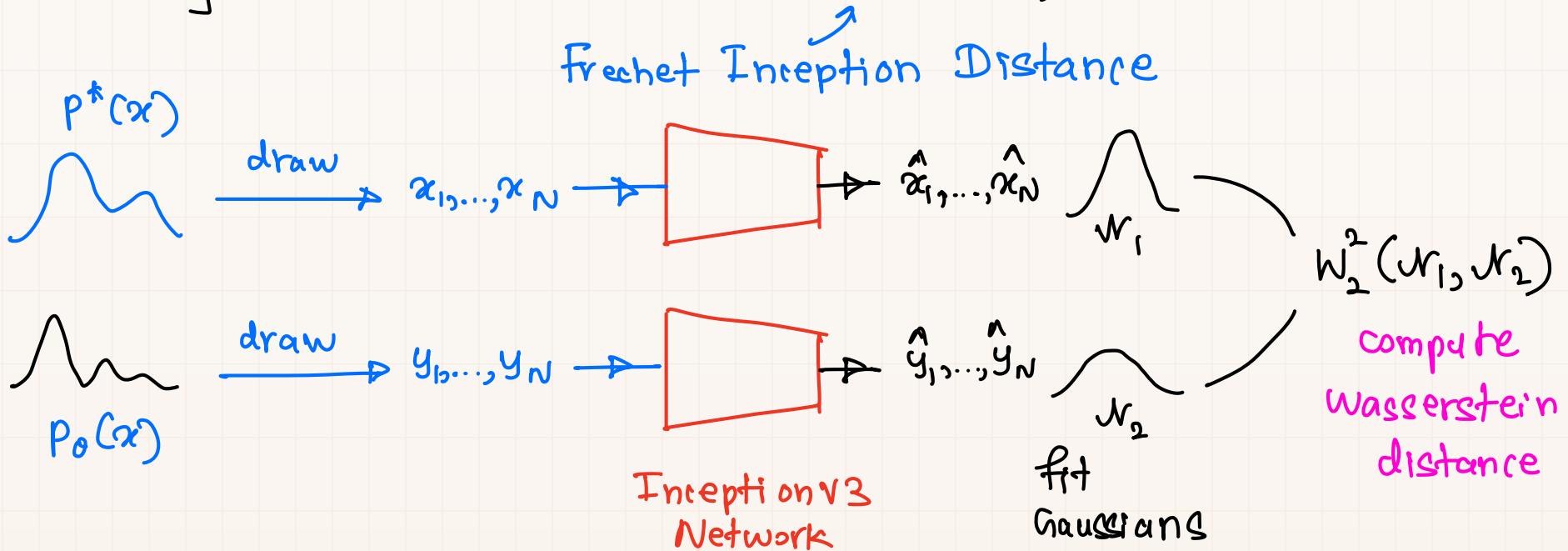
```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Notice the similarity
to the Langevin Monte-Carlo

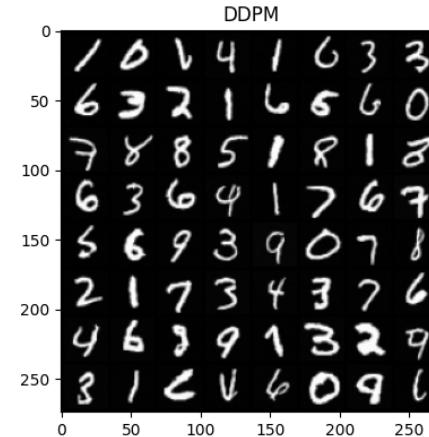
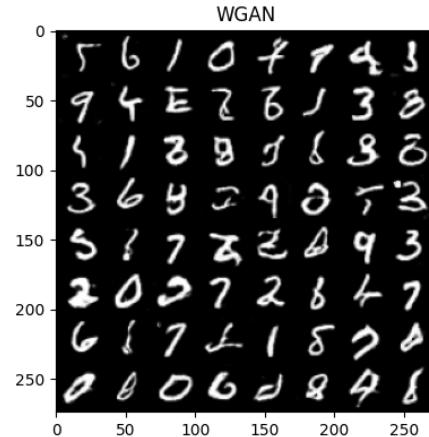
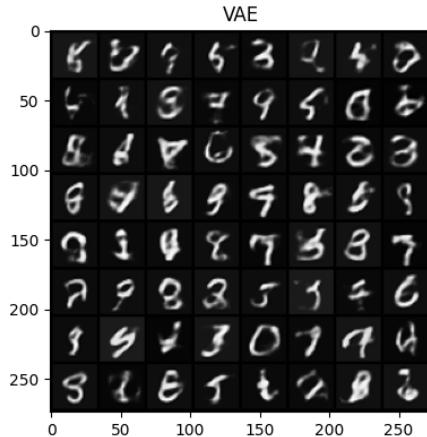
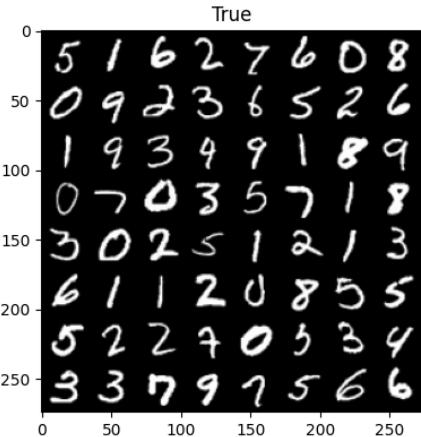
4

- Currently the state-of-the-art in generating images.
- Sampling is very slow
- Active Research Area.

Evaluating Generative Models (FID Score)



Some Results



FID

16.3

FID

12.4

FID

13.4

Code

main 1 Branch 0 Tags

Go to file Add file Code

Laknath1996 Update 4c99d02 · 13 hours ago 13 Commits

- ddpm Update 13 hours ago
- vae Update vae 5 days ago
- wgan Add ddpm scripts 2 weeks ago
- .gitignore Update ddpm 2 weeks ago
- README.md Update 2 weeks ago
- analysis.ipynb Update 13 hours ago

README

Description

Course project for EN.553.741 Machine Learning II at JHU. This repo contains code for implementing Wasserstein GAN, variational autoencoder and denoising diffusion probabilistic models from scratch.

Training metrics

About

Course project for EN.553.741 Machine Learning II at JHU. Implements Wasserstein GAN, variational autoencoder and denoising diffusion probabilistic models.

generative-adversarial-network
generative-model

Readme Activity 0 stars 2 watching 0 forks

Releases

No releases published [Create a new release](#)

Packages

No packages published [Publish your first package](#)



link to
repo

THANK YOU!

