Topics:

• Variational Autoencoders

CS 4803-DL / 7643-A ZSOLT KIRA

• Assignment 3

- Due July 13th 11:59pm EST
- Projects
 - Project proposal due July 26th

- Next Meta office hours today 3pm ET on bias/fairness
 - NOT recorded!

W9: July 7

Generative Models (Part I): Generative Adversarial Networks READING: LLaVA-OneVision

W9: July 9 Generative Models (Part II): Diffusion Models PS3/HW3 due July 13th 11:59pm (grace period July 15th)

W10: July 14 Variational Autoencoders (VAEs)

W10: July 16 Open topic! (please suggest). Otherwise default is Object Detection and Segmentation

 W11: July 21 Fairness/Bias Discussion, open topics not covered, Wrap-up
 Final Project Due July 26th 11:59pm (grace period July 28th)

Back to Generative Models





Supervised







Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





Variational Autoencoders (VAEs)





Goodfellow, NeurIPS 2016 Tutorial: Generative Adversarial Networks





Comparison





Low dimensional embedding

Linear layers with reduced dimension or Conv-2d layers with stride

Linear layers with increasing dimension or Conv-2d layers with bilinear upsampling





What is this? Hidden/Latent variables Factors of variation that produce an image: (digit, orientation, scale, etc.)



$$P(X) = \int P(X|Z;\theta)P(Z)dZ$$

We cannot maximize this likelihood due to the integral
Instead we maximize a variational *lower bound* (VLB) that we *can* compute

Kingma & Welling, Auto-Encoding Variational Bayes





- We can combine the probabilistic view, sampling, autoencoders, and approximate optimization
- Just as before, sample Z from simpler distribution
- We can also output parameters of a probability distribution!
 - **Example**: μ, σ of Gaussian distribution
 - For multi-dimensional version output diagonal covariance
- How can we maximize $P(X) = \int P(X|Z;\theta)P(Z)dZ$







 We can combine the probabilistic view, sampling, autoencoders, and approximate optimization



- Given an image, estimate Z
- Again, output parameters of a distribution





• We can tie the encoder and decoder together into a probabilistic autoencoder

- Given data (X), estimate μ_z, σ_z and sample from $N(\mu_z, \sigma_z)$
- Given Z, estimate μ_x , σ_x and sample from $N(\mu_x, \sigma_x)$







How can we optimize the parameters of the two networks?

Now equipped with our encoder and decoder networks, let's work out the (log) data likelihood:

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$



From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeur g



$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Logarithms}) \end{split}$$

From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeur g



Georg

Aside: KL Divergence (distance measure for distributions), always >= 0

$$KL(a||b) = H_c(a,b) - H(a) = \sum a(x)\log a(x) - \sum a(x)\log b(x)$$

Definition of Expectation

$$\mathbb{E}[f] = \mathbb{E}_{x \sim q}[f(x)] = \sum_{x \in \Omega} q(x) f(x)$$

$$KL(a||b) = E[\log a(x)] - E[\log b(x)] = E[\log \frac{a(x)}{b(x)}]$$





Maximizing Likelihood

Georg

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)})\right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Bayes' Rule})$$

$$= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})}\right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)}\right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})}\right] \quad (\text{Logarithms})$$

$$= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$\stackrel{\clubsuit}{=} \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$\stackrel{\clubsuit}{=} \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$\stackrel{\clubsuit}{=} \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

$$\stackrel{\clubsuit}{=} \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z)\right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))$$

From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeur g

Maximizing Likelihood

Georg

$$\begin{split} \log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} \left[\log p_{\theta}(x^{(i)}) \right] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_{z} \left[\log \frac{p_{\theta}(x^{(i)} \mid z)p_{\theta}(z)}{p_{\theta}(z \mid x^{(i)})} \frac{q_{\phi}(z \mid x^{(i)})}{q_{\phi}(z \mid x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_{z} \left[\log \frac{q_{\phi}(z \mid x^{(i)})}{p_{\theta}(z \mid x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))}_{>0} \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))}_{>0} \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid |p_{\theta}(z \mid x^{(i)}))}_{>0} \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))}_{>0} \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}))}_{>0} \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right]_{z} + \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(z \mid x^{(i)}) \mid p_{\theta}(z \mid x^{(i)}) \right] \right] \\ &= \underbrace{\mathbf{E}_{z} \left[\log p_{\theta}(x^{(i)} \mid z) \right]_{z} + \underbrace{\mathbf{E}_{z} \left[\log p_{\theta$$

From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeur g

0

GE

Maximizing Likelihood

Putting it all together: maximizing the likelihood lower bound



From CS231n, Fei-Fei Li, Justin Johnson, Serena Yeur g

0



Putting it all together: maximizing the likelihood lower bound



0



- Problem with respect to the VLB: updating ϕ $\mathcal{L}_{\text{VAE}} = \mathbb{E}_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \left[\log \frac{p_{\theta}(\boldsymbol{z}, \boldsymbol{x})}{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})} \right]$
 - $= -D_{ ext{KL}}(q_{\phi}(oldsymbol{z}|oldsymbol{x})||p_{ heta}(oldsymbol{z})) + \mathbb{E}_{q_{\phi}(oldsymbol{z}|oldsymbol{x})}[\log p_{ heta}(oldsymbol{x}|oldsymbol{z})]$
- $Z \sim Q(Z|X; \phi)$: need to differentiate through the sampling process w.r.t ϕ (encoder is probabilistic)



From: Tutorial on Variational Autoencoders <u>https://arxiv.org/abs/1606.05908</u>

From: http://gokererdogan.github.io/2016/07/01/reparameterization-trick/





- Solution: make the randomness independent of encoder output, making the encoder deterministic
- Gaussian distribution example:
 - Previously: encoder output = random variable z~N(μ, σ)
 - Now encoder output = distribution parameter [μ, σ]
 - $z = \mu + \epsilon * \sigma, \epsilon \sim N(0,1)$



From: Tutorial on Variational Autoencoders <u>https://arxiv.org/abs/1606.05908</u>

From: http://gokererdogan.github.io/2016/07/01/reparameterization-trick/







*Z*₂

Kingma & Welling, Auto-Encoding Variational Bayes





 Z_1

- Variational Autoencoders (VAEs) provide a principled way to perform approximate maximum likelihood optimization
 - Requires some assumptions (e.g. Gaussian distributions)
- Samples are often not as competitive as diffusion models or GANs
- Latent features (learned in an unsupervised way!) often good for downstream tasks:
 - Example: World models for reinforcement learning (Ha et al., 2018)



Summarv



De-noising Auto-encoder



Slide by Hung-yi Lee

Vincent, Pascal, et al. "Extracting and composing robust features with denoising autoencoders." *ICML*, 2008.



Discrete Representation

• Vector Quantized Variational Auto-encoder (VQVAE)



https://arxiv.org/abs/1711.00937

Slide by Hung-yi Lee



- Variational Autoencoders (VAEs) provide a principled way to perform approximate maximum likelihood optimization
 - Requires some assumptions (e.g. Gaussian distributions)
- Samples are often not as competitive as GANs
- Latent features (learned in an unsupervised way!) often good for downstream tasks:
 - Example: World models for reinforcement learning (Ha et al., 2018)







Q. Which ones are VAEs good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations			
Fast sampling			
High quality samples			



VAEs are bad at generating high quality samples

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations			
Fast sampling			
High quality samples	×		

Q. Which ones are GANs good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations			
Fast sampling			
High quality samples	×		

GANs suffer from mode collapse

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations		×	
Fast sampling		\checkmark	
High quality samples	×		

Q. Which ones are Diffusion models good at?

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations		×	
Fast sampling	\checkmark	\checkmark	
High quality samples	×		

Diffusion models are bad at sampling fast.

	Autoregressive (VAEs)	GANs	Diffusion
Mode coverage / diversity of generations		×	
Fast sampling	\checkmark	\checkmark	×
High quality samples	×		

Several ways to learn generative models via deep learning

Generative Adversarial Networks (GANs):

- Pro: Amazing results across many image modalities
- Con: Unstable/difficult training process, computationally heavy for good results
- Con: Limited success for discrete distributions (language)
- Con: Hard to evaluate (implicit model)

• Variational Autoencoders:

- Pro: Principled mathematical formulation
- Pro: Results in disentangled latent representations
- Con: Approximation inference, results in somewhat lower quality reconstructions

Diffusion Models

- Pro: Great results and diversity!
- Con: Slow generation (though lots of tricks to address)

Ha & Schmidhuber, World Models, 2018





Comparison





Denoising Diffusion Probabilistic Models (DDPMs)

And Conditional Diffusion Models

TEXT DESCRIPTION

An astronaut Teddy bears A bowl of soup

riding a horse lounging in a tropical resort in space playing basketball with cats in space

in a photorealistic style in the style of Andy Warhol as a pencil drawing DALL·E 2





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

Landscape Highlights of Diffusion Models (Nov 2022)

- Diffusion probabilistic models (Sohl-Dickstein et al., 2015)
- Noise-conditioned score network (NCSN; Yang & Ermon, 2019)
- Denoising diffusion probabilistic models (DDPM; <u>Ho et al. 2020</u>)
 - Classifier-guided conditional generation (Dhariwal and Nichole, 2021)
 - Classifier-free Diffusion Guidance (Ho and Salimans, 2022)
 - Latent-space Diffusion (StableDiffusion; Rombach and Blattmann et al., 2022)
 - Planning with Diffusion for Flexible Behavior Synthesis (Diffuser; Janner et al., 2022)
- new applications
- DreamFusion: Text-to-3D using 2D Diffusion (Poole and Jain et al., 2022)
 - Make-A-Video: Text-to-Video Generation without Text-Video Data (Singer et al., 2022)

basic principles

conditional & high-res image generation

How to make a new generative model

- **Setting:** Given unlabeled dataset of data, I want to learn to sample from P(x)
- Define the generative process
- Parameterize it
- Maximum likelihood (often -> KL-divergence)
- Approximations
- Optimize parameters!
- Add conditioning, e.g. text

Landscape Highlights of Diffusion Models (Nov 2022)

- basic _ Noise
- conditional & high-res image generation
- new applications

- Diffusion probabilistic models (<u>Sohl-Dickstein et al., 2015</u>)
- Noise-conditioned score network (NCSN; Yang & Ermon, 2019)
 - Denoising diffusion probabilistic models (DDPM; <u>Ho et al. 2020</u>)
 - Classifier-guided conditional generation (Dhariwal and Nichole, 2021)
 - Classifier-free Diffusion Guidance (Ho and Salimans, 2022)
 - Latent-space Diffusion (StableDiffusion; Rombach and Blattmann et al., 2022)
 - Planning with Diffusion for Flexible Behavior Synthesis (Diffuser; Janner et al., 2022)
- DreamFusion: Text-to-3D using 2D Diffusion (Poole and Jain et al., 2022)
 - Make-A-Video: Text-to-Video Generation without Text-Video Data (Singer et al., 2022)

image from dataset





image from dataset

The "forward diffusion" process: add Gaussian noise each step



image from dataset

The "forward diffusion" process: add Gaussian noise each step

noise $\mathcal{N}(0, I)$



 \bullet \bullet \bullet



image from dataset

The "forward diffusion" process: add Gaussian noise each step

noise $\mathcal{N}(0, I)$



 $X_1 \longleftarrow$ x_0

• • •

 $\xrightarrow{} x_{T-1} \xrightarrow{} x_T$

 $\longrightarrow x_{T-1} \longleftarrow x_T$

The "denoising diffusion" process: generate an image from noise by *denoising* the gaussian noises Ties/inspiration form Annealed Imporantce Sampling in physics

Comparison











The **known** forward process

$$x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$$

The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t} \)x_{t-1}, \beta_t I)$ Conditional Gaussian

The **known** forward process
$$x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$$

 $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)
 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t} \)x_{t-1}, \beta_t I)$ Conditional Gaussian

Notation: A Gaussian distribution "for" x_t

The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I)$ Conditional Gaussian

 β_t is the *variance schedule* at the diffusion step t

The **known** forward process
$$x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$$

 $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I)$ Conditional Gaussian



https://www.youtube.com/watch?v=HoKDTa5jHvg&t=517s

The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I)$ Conditional Gaussian

 β_t is the *variance schedule* at the diffusion step t

 $0 < \beta_1 < \beta_2 < \cdots < \beta_T < 1$, typical value range [0.0001, 0.02], with T = 1000



The **known** forward process $x_0 \longrightarrow x_1 \longrightarrow \cdots \longrightarrow x_T$ $q(x_{1:T}|x_0) = \prod_{t=1}^T q(x_t|x_{t-1})$ Probability Chain Rule (Markov Chain)

 $q(x_t|x_{t-1}) = \mathcal{N}(x_t; (\sqrt{1-\beta_t})x_{t-1}, \beta_t I)$ Conditional Gaussian

Nice property: samples from an *arbitrary forward step* are also Gaussian-distributed! $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\overline{\alpha}_t}x_0, (1 - \overline{\alpha}_t)I)$

, where $\alpha_t = (1 - \beta_t)$, $\overline{\alpha}_t = \prod_{s=1}^t \alpha_s$



$$\begin{aligned} & \underset{\alpha_{t}}{\text{https://www.youtube.com/watch?v=HokDTa5jHvg&t=517s} \\ & \underset{\alpha_{t}}{\alpha_{t}} = \prod_{s=1}^{T} a_{s} & q(x_{t}|x_{t-1}) = \mathcal{N}(x_{t}, \sqrt{1-\beta_{t}}x_{t-1}, \beta_{t}I) \\ & = \sqrt{1-\beta_{t}}x_{t-1} + \sqrt{\beta_{t}}\epsilon \\ & = \sqrt{\alpha_{t}}x_{t-1} + \sqrt{1-\alpha_{t}}\epsilon \\ & = \sqrt{\alpha_{t}}\alpha_{t-1}x_{t-2} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\epsilon \\ & = \sqrt{\alpha_{t}}\alpha_{t-1}x_{t-2} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\epsilon \\ & = \sqrt{\alpha_{t}}\alpha_{t-1}\alpha_{t-2}x_{t-3} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\alpha_{t-2}\epsilon \\ & = \sqrt{\alpha_{t}}\alpha_{t-1}\dots\alpha_{1}\alpha_{0}x_{0} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\dots\alpha_{1}\alpha_{0}\epsilon \\ & = \sqrt{\alpha_{t}}\alpha_{t-1}\dots\alpha_{1}\alpha_{0}x_{0} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\dots\alpha_{1}\alpha_{0}\epsilon \\ & = \sqrt{\alpha_{t}}x_{0} + \sqrt{1-\alpha_{t}}\alpha_{t-1}\dots\alpha_{1}\alpha_{0}\epsilon \\ & = \sqrt{\alpha_{t}}x_{0} + \sqrt{1-\alpha_{t}}\epsilon \\ & = \sqrt{\alpha_{t}}x_{0} +$$

Nice property: samples from an *arbitrary forward step* are also Gaussian-distributed! $q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\overline{\alpha}_t}x_0, (1 - \overline{\alpha}_t)I)$

Gaussian reparameterization trick:

$$z = \mu + \epsilon * \sigma, \epsilon \sim N(0,1)$$

Intuition: We know all distributions in forward process, and can in fact directly compute for any t based on X₀

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \qquad \epsilon \sim \mathcal{N}(0, I)$$

(square root appears because reparameterization trick has just σ)

The Diffusion and Denoising Process

