Lecture 19: Generative AI Part 2 GANs & Diffusion

Ranjay Krishna

Lecture 19 - 1



Administrative

- A5 is out. It is the last assignment.
- A5 deadline June 9th 11:59pm

• Almost done with the course :(

2

June 05, 2025

Lecture 19 -

Ranjay Krishna

Administrative

Final project details:

- **Poster session** June 9th, 10:30AM 12:20PM at Microsoft Atrium (CSE1)
 - Upload PDFs here by tomorrow morning to have us print them
 - Upload PDFs to <u>Gradescope</u> by Monday morning for grading purposes
- Final reports due June 9th

Ranjay Krishna

Lecture 19 - 3



Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

4

Ranjay Krishna

Ffjord

Lecture 19 - 4

Generative AI so far: Autoregressive models

Lecture 19 - 5

Generate image pixels starting from corner

Dependency on previous pixels modeled using an RNN (LSTM)

Very slow during both training and testing; N x N image requires 2N-1 sequential steps!



5

[van der Oord et al. 2016]

June 05, 2025

Ranjay Krishna

Variational Autoencoders: Intractability

Data likelihood:
$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Another idea: $p_{\theta}(x) = \frac{p_{\theta}(x|z)p_{\theta}(z)}{p_{\theta}(z|x)}$ **x**: 28x28 image = 784-dim vector
z: 20-dim vector
Encoder Network
 $q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$ **Decoder Network**
 $q_{\phi}(z|x) = N(\mu_{z|x}, \Sigma_{z|x})$ $p_{\theta}(x|z) = N(\mu_{x|z}, \Sigma_{x|z})$
 $\mu_{z|x}: 20$ $\Sigma_{z|x}: 20$ $\mu_{x|z}: 768$ $\Sigma_{x|z}: 768$
Linear(400->20) Linear(400->20) Linear(400->768) Linear(400->768) Linear(400->768) Linear(400->768)

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Ranjay Krishna

Lecture 19 - 6

Today: implicit density models



Lecture 19 - 7



Generative Adversarial Networks (GANs)



Lecture 19 - 8



All the models together

Autoregressive models define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent **z**: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

Ranjay Krishna

Lecture 19 - 9

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent z:

$$p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample?

Ranjay Krishna

Lecture 19 - 10

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent z:

$$p_{ heta}(x) = \int p_{ heta}(z) p_{ heta}(x|z) dz$$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

What if we give up on explicitly modeling density, and just want ability to sample? GANs: not modeling any explicit density function!

Ranjay Krishna

Lecture 19 - 11



Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Ranjay Krishna

Ffjord

Lecture 19 - 12

Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .



Lecture 19 - 13

Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior p(z) (e.g. assume z is a multivariate gaussian). Sample z ~ p(z) and pass to a Generator Network x = G(z)

Then x is a sample from the Generator distribution p_{G} . We just need to make sure $p_{G} = p_{data}$!



Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior p(z) (e.g. assume z is a multivariate gaussian). Sample z ~ p(z) and pass to a Generator Network x = G(z)

Then x is a sample from the Generator distribution p_G . We just need to make sure $p_G = p_{data}!$ Generator Network

Train Generator Network G to convert z into fake data x sampled from p_G

Ranjay Krishna

Lecture 19 - 15

Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior p(z) (e.g. assume z is a multivariate gaussian). Sample z ~ p(z) and pass to a Generator Network x = G(z)

Then x is a sample from the Generator distribution p_{G} . We just need to make sure $p_{G} = p_{data}$!



Ranjay Krishna

Lecture 19 - 16

Setup: Assume we have data x_i drawn from distribution $p_{data}(x)$. Want to sample from p_{data} .

Idea: Introduce a latent variable z with simple prior p(z) (e.g. assume z is a multivariate gaussian). Sample z ~ p(z) and pass to a Generator Network x = G(z)

Then x is a sample from the Generator distribution p_{G} . We just need to make sure $p_{G} = p_{data}$!



Ranjay Krishna

Lecture 19 - 17

Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim p_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$



Jointly train generator G and discriminator D with a minimax game



Jointly train generator G and discriminator D with a minimax game



Jointly train generator G and discriminator D with a minimax game



Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim p_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} \left(\boldsymbol{G}(\boldsymbol{z}) \right) \right) \right] \right)$$





Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim p_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$
$$= \min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \boldsymbol{V}(\boldsymbol{G}, \boldsymbol{D})$$

Ranjay Krishna

Lecture 19 - 23

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim p_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} \left(\boldsymbol{G}(\boldsymbol{z}) \right) \right) \right] \right)$$

$$= \min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \boldsymbol{V}(\boldsymbol{G}, \boldsymbol{D}) \qquad \text{For t in 1, ... T:}$$

$$1. (\text{Update D}) \boldsymbol{D} = \boldsymbol{D} + \alpha_{\boldsymbol{D}} \frac{\partial \boldsymbol{V}}{\partial \boldsymbol{D}}$$

$$2. (\text{Update G}) \boldsymbol{G} = \boldsymbol{G} - \alpha_{\boldsymbol{G}} \frac{\partial \boldsymbol{V}}{\partial \boldsymbol{G}}$$

June 05, 2025

Lecture 19 - 24

Ranjay Krishna

Jointly train generator G and discriminator D with a minimax game

Train G and D using alternating gradient updates

$$\min_{G} \max_{D} \left(E_{x \sim p_{data}} [\log D(x)] + E_{z \sim p(z)} \left[\log \left(1 - D(G(z)) \right) \right] \right)$$

$$= \min_{G} \max_{D} V(G, D)$$
For t in 1, ... T:
$$1. (Update D) D = D + \alpha_{D} \frac{\partial V}{\partial D}$$
We are not minimizing any overall loss! No training curves to look at!
$$2. (Update G) G = G - \alpha_{G} \frac{\partial V}{\partial G}$$

Ranjay Krishna

Lecture 19 - 25

Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim p_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim p(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so D(G(z)) close to 0



June 05, 2025

Lecture 19 - 26

Ranjay Krishna

Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim \boldsymbol{p}_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim \boldsymbol{p}(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$

Lecture 19 - 27

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so D(G(z)) close to 0 **Problem**: Why is this a problem?

Ranjay Krishna



Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim \boldsymbol{p}_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim \boldsymbol{p}(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so D(G(z)) close to 0 **Problem**: Vanishing gradients for G How do we fix this?



June 05, 2025

Lecture 19 - 28

Ranjay Krishna

Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim \boldsymbol{p}_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim \boldsymbol{p}(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so D(G(z)) close to 0 **Problem:** Vanishing gradients for G **Solution:** Train G to minimize –log(D(G(z)), instead of log(1-D(G(z)). Then G gets strong gradients at start of training!

Ranjay Krishna



June 05, 2025

Lecture 19 - 29

Jointly train generator G and discriminator D with a minimax game

$$\min_{\boldsymbol{G}} \max_{\boldsymbol{D}} \left(E_{\boldsymbol{x} \sim \boldsymbol{p}_{data}} [\log \boldsymbol{D}(\boldsymbol{x})] + E_{\boldsymbol{z} \sim \boldsymbol{p}(\boldsymbol{z})} \left[\log \left(1 - \boldsymbol{D} (\boldsymbol{G}(\boldsymbol{z})) \right) \right] \right)$$

At start of training, generator is very bad and discriminator can easily tell apart real/fake, so D(G(z)) close to 0 **Problem:** Vanishing gradients for G **Solution:** Train G to minimize –log(D(G(z)), instead of log(1-D(G(z)). Then G gets strong gradients at start of training!

Ranjay Krishna



June 05, 2025

Lecture 19 - 30

Once trained, throw away the discriminator and use G to generate new images



Ranjay Krishna

Lecture 19 - 31

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generative Adversarial Nets

Generated samples



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Ranjay Krishna

Lecture 19 - 32

Ian Goodfellow et al., "Generative Adversarial Nets", NIPS 2014

Generative Adversarial Nets

Generated samples (CIFAR-10)



Nearest neighbor from training set

Figures copyright Ian Goodfellow et al., 2014. Reproduced with permission.

Ranjay Krishna

Lecture 19 - 33

Generative Adversarial Nets: Convolutional Architectures



Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Ranjay Krishna

Lecture 19 - 34

Generative Adversarial Nets: Convolutional Architectures

Generator is an upsampling network with fractionally-strided convolutions **Discriminator** is a convolutional network

Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

Radford et al, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks", ICLR 2016

Ranjay Krishna

Lecture 19 - 35

Generative Adversarial Nets: Convolutional Architectures

Samples from the model look much better!

Radford et al, ICLR 2016



Ranjay Krishna

Lecture 19 - 36
Generative Adversarial Nets: Convolutional Architectures

Interpolating between random points in laten space

Radford et al, ICLR 2016



Ranjay Krishna

Lecture 19 - 37

Radford et al, ICLR 2016

Samples from the model







Neutral man

Ranjay Krishna

Lecture 19 - 38

Neutral man

Radford et al, ICLR 2016

Samples from the model



Smiling woman Neutral woman

Ranjay Krishna

Lecture 19 - 39



Ranjay Krishna

Lecture 19 - 40

Glasses man

No glasses man

No glasses woman

Radford et al, ICLR 2016

Woman with glasses



Ranjay Krishna

Lecture 19 - 41

Since then: Explosion of GANs

"The GAN Zoo"

See also: <u>https://github.com/soumith/ganhacks</u> for tips and tricks for trainings GANs

- GAN Generative Adversarial Networks
- 3D-GAN Learning a Probabilistic Latent Space of Object Shapes via 3D Generative-Adversarial Modeling
- acGAN Face Aging With Conditional Generative Adversarial Networks
- AC-GAN Conditional Image Synthesis With Auxiliary Classifier GANs
- AdaGAN AdaGAN: Boosting Generative Models
- AEGAN Learning Inverse Mapping by Autoencoder based Generative Adversarial Nets
- AffGAN Amortised MAP Inference for Image Super-resolution
- · AL-CGAN Learning to Generate Images of Outdoor Scenes from Attributes and Semantic Layouts
- ALI Adversarially Learned Inference
- AM-GAN Generative Adversarial Nets with Labeled Data by Activation Maximization
- AnoGAN Unsupervised Anomaly Detection with Generative Adversarial Networks to Guide Marker Discovery
- ArtGAN ArtGAN: Artwork Synthesis with Conditional Categorial GANs
- b-GAN b-GAN: Unified Framework of Generative Adversarial Networks
- Bayesian GAN Deep and Hierarchical Implicit Models
- BEGAN BEGAN: Boundary Equilibrium Generative Adversarial Networks
- BiGAN Adversarial Feature Learning
- BS-GAN Boundary-Seeking Generative Adversarial Networks
- CGAN Conditional Generative Adversarial Nets
- CaloGAN CaloGAN: Simulating 3D High Energy Particle Showers in Multi-Layer Electromagnetic Calorimeters
 with Generative Adversarial Networks
- CCGAN Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks
- CatGAN Unsupervised and Semi-supervised Learning with Categorical Generative Adversarial Networks
- CoGAN Coupled Generative Adversarial Networks

- Context-RNN-GAN Contextual RNN-GANs for Abstract Reasoning Diagram Generation
- C-RNN-GAN C-RNN-GAN: Continuous recurrent neural networks with adversarial training
- CS-GAN Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets
- CVAE-GAN CVAE-GAN: Fine-Grained Image Generation through Asymmetric Training
- CycleGAN Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks
- DTN Unsupervised Cross-Domain Image Generation
- DCGAN Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks
- DiscoGAN Learning to Discover Cross-Domain Relations with Generative Adversarial Networks
- DR-GAN Disentangled Representation Learning GAN for Pose-Invariant Face Recognition
- DualGAN DualGAN: Unsupervised Dual Learning for Image-to-Image Translation
- EBGAN Energy-based Generative Adversarial Network
- f-GAN f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization
- FF-GAN Towards Large-Pose Face Frontalization in the Wild
- GAWWN Learning What and Where to Draw
- GeneGAN GeneGAN: Learning Object Transfiguration and Attribute Subspace from Unpaired Data
- Geometric GAN Geometric GAN
- GoGAN Gang of GANs: Generative Adversarial Networks with Maximum Margin Ranking
- GP-GAN GP-GAN: Towards Realistic High-Resolution Image Blending
- IAN Neural Photo Editing with Introspective Adversarial Networks
- iGAN Generative Visual Manipulation on the Natural Image Manifold
- IcGAN Invertible Conditional GANs for image editing
- ID-CGAN Image De-raining Using a Conditional Generative Adversarial Network
- Improved GAN Improved Techniques for Training GANs
- InfoGAN InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets
- LAGAN Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics
 Synthesis
- LAPGAN Deep Generative Image Models using a Laplacian Pyramid of Adversarial Networks

https://github.com/hindupuravinash/the-gan-zoo

June 05, 2025

Ranjay Krishna

Lecture 19 - 42

GAN improvements: better loss functions





LSGAN, Zhu 2017. Wasserstein GAN, Arjovsky 2017.

Improved Wasserstein GAN, Gulrajani 2017.

Ranjay Krishna

Lecture 19 - 43

GAN improvements: higher resolution

256 x 256 bedrooms



1024 x 1024 faces



Progressive GAN, Karras 2018.

Ranjay Krishna

Lecture 19 - 44

GAN transformations

Source->Target domain transfer





zebra → horse



apple \rightarrow orange









→ summer Yosemite



CycleGAN. Zhu et al. 2017.



Pix2pix. Isola 2017. Many examples at https://phillipi.github.io/pix2pix/

Ranjay Krishna

Lecture 19 - 45

BigGAN: 512x512 images



Brock et al., 2019

Ranjay Krishna

Lecture 19 - 46

GANs with self-attention mechanism



Zhang et al, "Self-Attention Generative Adversarial Networks", ICML 2019

Ranjay Krishna

Lecture 19 - 47

Controlled generation with GANs



Semantic Manipulation Using Segmentation Map

Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Ranjay Krishna

Lecture 19 - 48

Controlled generation with GANs



Semantic Manipulation Using Segmentation Map

Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Ranjay Krishna

Lecture 19 - 49

Controlled generation with GANs



Park et al, "Semantic Image Synthesis with Spatially-Adaptive Normalization", CVPR 2019

Ranjay Krishna

Lecture 19 - 50

Conditional GANs: StyleGAN

Y is text that describes the image you want to generate



Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Ranjay Krishna

Lecture 19 - 51

Conditional GANs: StyleGAN

Batch Normalization

$\mu_j = \frac{1}{N} \sum_{i=1}^{N} x_{i,j}$ $\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{L-1} (x_{i,j} - \mu_j)^2$ $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$ $y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$

Learn a separate scale and shift for each different label y **Conditional Batch Normalization**

June 05, 2025

 $\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$ $\sigma_j^2 = \frac{1}{N} \sum_{i=1}^{N} (x_{i,j} - \mu_j)^2$ $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$ $y_{i,j} = \frac{\gamma_j^y}{\gamma_j} \hat{x}_{i,j} + \beta_j^y$

Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Ranjay Krishna

Lecture 19 - 52

Conditional GANs: StyleGAN



Karras et al, "Analyzing and Improving the Image Quality of StyleGAN", CVPR 2020

Ranjay Krishna

Lecture 19 - 53

Scene graphs to GANs

Specifying exactly what kind of image you want to generate.

The explicit structure in scene graphs provides better image generation for complex scenes.

Scene Graph





Figures copyright 2019. Reproduced with permission.

Johnson et al. Image Generation from Scene Graphs, CVPR 2019

Ranjay Krishna

Lecture 19 - 54

HYPE: Human eYe Perceptual Evaluations



Lecture 19 - 55

Zhou, Gordon, Krishna et al. HYPE: Human eYe Perceptual Evaluations, NeurIPS 2019

Figures copyright 2019. Reproduced with permission.

June 05, 2025

Ranjay Krishna

Summary: GANs

Pros:

- Beautiful samples, was state-of-the-art until diffusion models!

Cons:

- Trickier / more unstable to train
- Can't solve inference queries such as p(x), p(z|x)

Active areas of research:

- Better loss functions, more stable training (Wasserstein GAN, LSGAN, many others)
- Conditional GANs, GANs for all kinds of applications

Ranjay Krishna

Lecture 19 - 56

Diffusion models







Diffusion Models are outperforming GANs



Dhariwal & Nichol. "Diffusion Models Beat GANs on Image Synthesis", OpenAI 2021



Ho et al. "Cascaded Diffusion Models for High Fidelity Image Generation", Google 2021

Ranjay Krishna

Lecture 19 - 58

Text-to-Image (T2I) Generation

Dall-E2

"a teddy bear on a skateboard in times square"



Ramesh et al. "Hierarchical Text-Conditional Image Generation with CLIP Latents" 2022

Imagen

"A group of teddy bears in suit in a corporate office celebrating the birthday of their friend. There is a pizza cake on the desk."



Saharia et al. "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding" 2022

Ranjay Krishna

Lecture 19 - 59

Text-to-Image (T2I) Generation

Stable Diffusion



Mega thread on Twitter/X about Stable Diffusion

Rombach et al. "High-Resolution Image Synthesis with Latent Diffusion Models" 2022

Ranjay Krishna

Lecture 19 - 60

Application of diffusion: Image Super-resolution

Irish Setter

Saharia et al., Image Super-Resolution via Iterative Refinement, ICCV 2021

Gif on this slide is not displayed in pdf

June 05, 2025

Ranjay Krishna

Lecture 19 - 61

But what is a diffusion model?







So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent **z**: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

GANs give up on explicitly modeling density and just learns to sample "real" data

Ranjay Krishna

Lecture 19 - 63

So far...

Autoregressive define tractable density function, optimize likelihood of training data:

$$p_{\theta}(x) = \prod_{i=1}^{n} p_{\theta}(x_i | x_1, ..., x_{i-1})$$

VAEs define intractable density function with latent **z**: $p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$

Cannot optimize directly, derive and optimize lower bound on likelihood instead

GANs give up on explicitly modeling density and just learns to sample "real" data

All these methods generate data in one forward step! Why this is hard?

Ranjay Krishna

Lecture 19 - 64



Figure copyright and adapted from Ian Goodfellow, Tutorial on Generative Adversarial Networks, 2017.

Ranjay Krishna

Ffjord

Lecture 19 - 65

Recall VAEs

VAEs define intractable density function with latent **z**: $p_{\theta}(x) = \int p_{\theta}(z) p_{\theta}(x|z) dz$

Cannot optimize directly, derive and optimize lower bound on likelihood instead:



First loss for the encoder



Second loss for both decoder and encoder



VAEs for images look like this



- We learn 2 networks, one to encode and one to decode
- We ensure that z is similar to a unit normal noise
- To sample new images, we can sample from the unit normal and decode in 1 step

Ranjay Krishna

Lecture 19 - 69

Markovian Hierarchical VAEs



- We learn 2T networks, one to encode and one to decode
- We ensure that z_{τ} is similar to a unit normal noise
- To sample new images, we can sample from the unit normal and decode in T step

Ranjay Krishna

Lecture 19 - 70

Markovian Hierarchical VAEs - same derivation

$$\log p(x) = \mathbb{E}_{z_{1:T} \sim q_{\phi}(z_{1:T}|x)} [\log p_{\theta}(x^{(i)})] \\ = \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{p_{\theta}(z_{1:T}|x)}]$$

 $p_{\theta}(x)$ is independent of $z_{1:T}$

Bayes rule

Ranjay Krishna

Lecture 19 - 71

Markovian Hierarchical VAEs - same derivation

$$\begin{split} \log p(x) &= \mathbb{E}_{z_{1:T} \sim q_{\phi}(z_{1:T}|x)}[\log p_{\theta}(x^{(i)})] & p_{\theta}(x) \text{ is independent of } z_{1:T} \\ &= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{p_{\theta}(z_{1:T}|x)}] & \text{Bayes rule} \\ &= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{p_{\theta}(z_{1:T}|x)} \frac{q_{\phi}(z_{1:T}|x)}{q_{\phi}(z_{1:T}|x)}] & \text{Multiplying by a constant} \\ &= \mathbb{E}_{z_{1:T}}[\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T}|x)}] + \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T})}{p_{\theta}(z_{1:T}|x)}] \end{split}$$

Ranjay Krishna

]

Lecture 19 - 72
Markovian Hierarchical VAEs - same derivation

$$\log p(x) = \mathbb{E}_{z_{1:T} \sim q_{\phi}(z_{1:T}|x)}[\log p_{\theta}(x^{(i)})] \qquad p_{\theta}(x) \text{ is independent of } z_{1:T}$$

$$= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{p_{\theta}(z_{1:T}|x)}] \qquad \text{Bayes rule}$$

$$= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{p_{\theta}(z_{1:T}|x)} \frac{q_{\phi}(z_{1:T}|x)}{q_{\phi}(z_{1:T}|x)}] \qquad \text{Multiplying by a constant}$$

$$= \mathbb{E}_{z_{1:T}}[\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}] + \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T})}{p(z_{1:T}|x)}]$$
Reconstruction objective maximizes the likelihood of data $p_{\theta}(x|z)$
This KL term (between Gaussians for encoder and z prior)

Ranjay Krishna

Lecture 19 - 73

Keeping just the first two terms:

$$\log p(x) \ge = \mathbb{E}_{z_{1:T}}[\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}]$$



Lecture 19 - 74

Keeping just the first two terms:

$$\begin{split} \log p(x) &\geq = \mathbb{E}_{z_{1:T}} [\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} [\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \end{split}$$

Ranjay Krishna

Lecture 19 - 75

Keeping just the first two terms:

$$\begin{split} \log p(x) &\geq = \mathbb{E}_{z_{1:T}} [\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} [\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \end{split}$$

where the joint probability distribution is: $p(\boldsymbol{x}, \boldsymbol{z}_{1:T}) = p(\boldsymbol{z}_T)p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z}_1)\prod_{t=2}^T p_{\boldsymbol{\theta}}(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_t)$

This is very similar to the autoregressive model formula

Ranjay Krishna

Lecture 19 - 76

Keeping just the first two terms:

$$\begin{split} \log p(x) &\geq = \mathbb{E}_{z_{1:T}} [\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} [\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \end{split}$$

where the joint probability distribution is: $p(\boldsymbol{x}, \boldsymbol{z}_{1:T}) = p(\boldsymbol{z}_T)p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_t)$ And the encoder posterior is: $q_{\boldsymbol{\phi}}(\boldsymbol{z}_{1:T} \mid \boldsymbol{x}) = q_{\boldsymbol{\phi}}(\boldsymbol{z}_1 \mid \boldsymbol{x}) \prod_{t=2}^T q_{\boldsymbol{\phi}}(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})$

Ranjay Krishna

Lecture 19 - 77

Keeping just the first two terms:

$$\begin{split} \log p(x) &\geq = \mathbb{E}_{z_{1:T}}[\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}}[\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})}] \\ &= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}}[\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \end{split}$$
 Why is this a hard objective to train?

where the joint probability distribution is: $p(\boldsymbol{x}, \boldsymbol{z}_{1:T}) = p(\boldsymbol{z}_T)p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_t)$ And the encoder posterior is: $q_{\boldsymbol{\phi}}(\boldsymbol{z}_{1:T} \mid \boldsymbol{x}) = q_{\boldsymbol{\phi}}(\boldsymbol{z}_1 \mid \boldsymbol{x}) \prod_{t=2}^T q_{\boldsymbol{\phi}}(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})$

Ranjay Krishna

Lecture 19 - 78

Keeping just the first two terms:

$$\begin{split} \log p(x) &\geq = \mathbb{E}_{z_{1:T}} [\log p_{\theta}(x|z_{1:T})] - \mathbb{E}_{z_{1:T}} [\log \frac{q_{\phi}(z_{1:T}|x)}{p_{\theta}(z_{1:T})} \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x|z_{1:T})p_{\theta}(z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_{1:T}|x)}}] \\ &= \mathbb{E}_{z_{1:T}} [\log \frac{p_{\theta}(x,z_{1:T})}{q_{\phi}(z_$$

Why is this a hard objective to train?1. There are too many networks to learn2. The objective function is expensive!3. It collapses easily!

where the joint probability distribution is: $p(\boldsymbol{x}, \boldsymbol{z}_{1:T}) = p(\boldsymbol{z}_T)p_{\boldsymbol{\theta}}(\boldsymbol{x} \mid \boldsymbol{z}_1) \prod_{t=2}^T p_{\boldsymbol{\theta}}(\boldsymbol{z}_{t-1} \mid \boldsymbol{z}_t)$ And the encoder posterior is: $q_{\boldsymbol{\phi}}(\boldsymbol{z}_{1:T} \mid \boldsymbol{x}) = q_{\boldsymbol{\phi}}(\boldsymbol{z}_1 \mid \boldsymbol{x}) \prod_{t=2}^T q_{\boldsymbol{\phi}}(\boldsymbol{z}_t \mid \boldsymbol{z}_{t-1})$

Ranjay Krishna

Lecture 19 - 79



- Why is this a hard objective to train?
- 1. There are too many networks to learn.
- 2. The objective function is expensive!
- 3. It collapses easily!

Ranjay Krishna

Lecture 19 - 80



Diffusion models are a special case With a more interpretable, simpler objective.

Ranjay Krishna

Lecture 19 - 81

1. The latent dimension size is exactly equal to the data dimension



Ranjay Krishna

Lecture 19 - 82

- 1. The latent dimension size is exactly equal to the data dimension
- 2. The encoders are pre-defined and not learned.





Ranjay Krishna

Lecture 19 - 83

3. Encoders are designed as a linear Gaussian model conditioned on the time step: Add noise at every time step



Ranjay Krishna

Lecture 19 - 84

How the forward step was designed:

So, given x_{t-1} we can sample x_t using:

$$x_t \sim \sqrt{\alpha_t} x_{t-1} + (1 - \alpha_t) \epsilon$$

where $\epsilon \sim \mathcal{N}(x; 0, I)$



Lecture 19 - 85



$$egin{aligned} oldsymbol{x}_t &= \sqrt{ar{lpha}_t}oldsymbol{x}_0 + \sqrt{1-ar{lpha}_t}oldsymbol{\epsilon}_0 \ &\sim \mathcal{N}(oldsymbol{x}_t;\sqrt{ar{lpha}_t}oldsymbol{x}_0,(1-ar{lpha}_t)\mathbf{I}) \end{aligned}$$

 \mathbf{x}_{t} is now a Gaussian characterized by \mathbf{x}_{0}

where
$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

Ranjay Krishna

Lecture 19 - 86

4. The Gaussian parameters vary over time in such a way that the distribution of the latent at final step T is a standard Gaussian



Ranjay Krishna

Lecture 19 - 87

Terminology: Forward and backward process

Forward diffusion process (fixed)



Reverse denoising process (generative)

Note: reverse or backward here doesn't mean the same thing as backpropagation

Ranjay Krishna

Lecture 19 - 88

The distribution perspective

Over time, as we add more noise sampled from a Gaussian distribution, it begins to look more like a unit normal



Ranjay Krishna

Lecture 19 - 89

Q. What do we have to learn to generate new samples from noise?



Ranjay Krishna

Lecture 19 - 90

Q. What do we have to learn to generate new samples from noise? A. We want to define a neural network to predict $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$



Ranjay Krishna

Lecture 19 - 91

Q. How should we train $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$?



Ranjay Krishna

Lecture 19 - 92

Q. How should we train $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$? A. We can get it to match $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0})$!



Ranjay Krishna

Lecture 19 - 93

Q. How should we train $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$? A. We can get it to match $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0})$!

Q. But why $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0)$ and not:



Ranjay Krishna

Lecture 19 - 94

Ok so our loss function is:

Q. How should we train $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$? A. We can get it to match $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0})$!

Minimize the distance between the two distributions:

$$rgmin_{oldsymbol{ heta}} \mathcal{D}_{ ext{KL}}(q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) \mid\mid p_{oldsymbol{ heta}}(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t))$$

Ranjay Krishna

Lecture 19 - 95

Ok so our loss function is:

Q. How should we train $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t})$? A. We can get it to match $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0})$!

Minimize the distance between the two distributions:

$$rgmin_{oldsymbol{ heta}} \mathcal{D}_{ ext{KL}}(q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) \mid\mid p_{oldsymbol{ heta}}(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t))$$

Problem: How do we estimate $q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0)$?

Ranjay Krishna

Lecture 19 - 96

The forward diffusion step

The distribution at step t is a Gaussian

$$q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}) = \mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 97

The forward diffusion step

The distribution at step t is a Gaussian

$$q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}) = \mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})$$

The mean defined by x_{t-1} : $\boldsymbol{\mu}_t(\boldsymbol{x}_t) = \sqrt{\alpha_t} \boldsymbol{x}_{t-1}$

 a_{t} is a predefined value for each step t

Ranjay Krishna

Lecture 19 - 98

The forward diffusion step

The distribution at step t is a Gaussian

$$q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}) = \mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})$$

The mean defined by $\mathsf{x}_{\mathsf{t} ext{-1}}$: $oldsymbol{\mu}_t(oldsymbol{x}_t) = \sqrt{lpha_t}oldsymbol{x}_{t-1}$

 a_{t} is a predefined value for each step t

The covariance is independent of x_{t-1} (an assumption)

$$oldsymbol{\Sigma}_t(oldsymbol{x}_t) = (1-lpha_t) \mathbf{I}_t$$

Ranjay Krishna

Lecture 19 - 99

How the forward step was designed:

The distribution at step t is a Gaussian

$$q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}) = \mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})$$

$$oldsymbol{\mu}_t(oldsymbol{x}_t) = \sqrt{lpha_t}oldsymbol{x}_{t-1}$$
 $oldsymbol{\Sigma}_t(oldsymbol{x}_t) = (1-lpha_t)oldsymbol{I}_t$

So, given x_{t-1} we can sample x_t using:

$$x_t \sim \sqrt{\alpha_t} x_{t-1} + (1 - \alpha_t) \epsilon$$

where $\epsilon \sim \mathcal{N}(x; 0, I)$

Ranjay Krishna

Lecture 19 - 100

 $oldsymbol{x}_t = \sqrt{lpha_t}oldsymbol{x}_{t-1} + \sqrt{1-lpha_t}oldsymbol{\epsilon}_{t-1}^*$



Lecture 19 - 101

$$oldsymbol{x}_t = \sqrt{lpha_t} oldsymbol{x}_{t-1} + \sqrt{1 - lpha_t} oldsymbol{\epsilon}_{t-1}^* \ = \sqrt{lpha_t} \left(\sqrt{lpha_{t-1}} oldsymbol{x}_{t-2} + \sqrt{1 - lpha_{t-1}} oldsymbol{\epsilon}_{t-2}^*
ight) + \sqrt{1 - lpha_t} oldsymbol{\epsilon}_{t-1}^*$$
 Substituting $oldsymbol{x}_{t-1}$

Ranjay Krishna

Lecture 19 - 102

Ranjay Krishna

Lecture 19 - 103

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}} \boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}} \left(\sqrt{\alpha_{t-1}} \boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^{*} \right) + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacktriangleleft \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}} \alpha_{t-1} \boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacksquare \qquad \text{Opening the parentheses} \end{split}$$

We can interpret this $\sqrt{1-\alpha_t} \epsilon_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0},(1-\alpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon^*_{t-2}$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$

Ranjay Krishna

Lecture 19 - 104

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}} \boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}} \left(\sqrt{\alpha_{t-1}} \boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^{*} \right) + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacktriangleleft \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}} \alpha_{t-1} \boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacksquare \qquad \text{Opening the parentheses} \end{split}$$

We can interpret this $\sqrt{1-lpha_t} m{\epsilon}^*_{t-1}$ as a sample from $\mathcal{N}(\mathbf{0},(1-lpha_t)\mathbf{I})$

We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon^*_{t-2}$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$

Notice that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^* + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ is the sum of two Gaussian samples

Ranjay Krishna

Lecture 19 - 105

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}} \boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}} \left(\sqrt{\alpha_{t-1}} \boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^{*} \right) + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacktriangleleft \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}} \alpha_{t-1} \boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}} \boldsymbol{\alpha}_{t-1} \boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacksquare \qquad \text{Opening the parentheses} \end{split}$$

We can interpret this $\sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$ We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$ Notice that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^* + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ is the sum of two Gaussian samples Using the property: $\mathcal{N}(x; 0, \sigma_1 I) + \mathcal{N}(x; 0, \sigma_2 I) = \mathcal{N}(x; 0, \sqrt{\sigma_1^2 + \sigma_2^2} I)$

Ranjay Krishna

Lecture 19 - 106

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}} \boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}} \left(\sqrt{\alpha_{t-1}} \boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}^{*} \right) + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacktriangleleft \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}} \alpha_{t-1} \boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}} \boldsymbol{\epsilon}_{t-1}^{*} \blacksquare \qquad \text{Opening the parentheses} \end{split}$$

We can interpret this $\sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (1 - \alpha_t)\mathbf{I})$ We can interpret this $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^*$ as a sample from $\mathcal{N}(\mathbf{0}, (\alpha_t - \alpha_t \alpha_{t-1})\mathbf{I})$ **Notice** that $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^* + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ is the sum of two Gaussian samples Using the property: $\mathcal{N}(x; 0, \sigma_1 I) + \mathcal{N}(x; 0, \sigma_2 I) = \mathcal{N}(x; 0, \sqrt{\sigma_1^2 + \sigma_2^2} I)$ We can rewrite $\sqrt{\alpha_t - \alpha_t \alpha_{t-1}} \epsilon_{t-2}^* + \sqrt{1 - \alpha_t} \epsilon_{t-1}^*$ as $\sqrt{\sqrt{\alpha_t - \alpha_t \alpha_{t-1}}^2} + \sqrt{1 - \alpha_t}^2 \epsilon_{t-2}$

Ranjay Krishna

Lecture 19 - 107

$$\begin{aligned} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}}\alpha_{t-1}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t} - \alpha_{t}}\alpha_{t-1}}^{2} + \sqrt{1 - \alpha_{t}}^{2}\boldsymbol{\epsilon}_{t-2} \quad \text{Sum of two Gaussians} \end{aligned}$$

Ranjay Krishna

Lecture 19 - 108
$$\begin{aligned} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1}}^{2} + \sqrt{1 - \alpha_{t}}^{2}}\boldsymbol{\epsilon}_{t-2} \quad \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1} + 1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-2} \quad \text{Squaring the terms} \end{aligned}$$

Ranjay Krishna

Lecture 19 - 109

$$\begin{aligned} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1-\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \quad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}}^{2} + \sqrt{1-\alpha_{t}}^{2}}\boldsymbol{\epsilon}_{t-2} \quad \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}+1-\alpha_{t}}\boldsymbol{\epsilon}_{t-2}} \quad \text{Squaring the terms} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1-\alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}} \quad \text{Simplifying} \end{aligned}$$

Ranjay Krishna

Lecture 19 - 110

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1-\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1-\alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}}^{2}} + \sqrt{1-\alpha_{t}}^{2}\boldsymbol{\epsilon}_{t-2} \qquad \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t}-\alpha_{t}\alpha_{t-1}} + 1-\alpha_{t}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Squaring the terms} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1-\alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Simplifying} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^{t}\alpha_{i}}\boldsymbol{x}_{0} + \sqrt{1-\prod_{i=1}^{t}\alpha_{i}}\boldsymbol{\epsilon}_{0} \qquad \text{Substituting till } \boldsymbol{x}_{0} \end{split}$$

Ranjay Krishna

Lecture 19 - 111

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}}\alpha_{t-1}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t} - \alpha_{t}}\alpha_{t-1}^{2}} + \sqrt{1 - \alpha_{t}}^{2}\boldsymbol{\epsilon}_{t-2} \qquad \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}}\alpha_{t-1} + 1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Squaring the terms} \\ &= \sqrt{\alpha_{t}}\alpha_{t-1}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t}}\alpha_{t-1} + 1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Simplifying} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^{t}}\alpha_{i}\boldsymbol{x}_{0} + \sqrt{1 - \prod_{i=1}^{t}}\alpha_{i}\boldsymbol{\epsilon}_{0} \qquad \text{Substituting till } \boldsymbol{x}_{0} \\ &= \sqrt{\alpha_{t}}\boldsymbol{x}_{0} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{0} \qquad \text{Let } \bar{\alpha}_{t} = \prod_{i=1}^{t}\alpha_{i} \end{split}$$

Ranjay Krishna

Lecture 19 - 112

$$\begin{split} \boldsymbol{x}_{t} &= \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \\ &= \sqrt{\alpha_{t}}\left(\sqrt{\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*}\right) + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Substituting } \boldsymbol{x}_{t-1} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2}^{*} + \sqrt{1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-1}^{*} \qquad \text{Opening the parentheses} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1}}^{2} + \sqrt{1 - \alpha_{t}}}\boldsymbol{\epsilon}_{t-2}^{2} \qquad \text{Sum of two Gaussians} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{\alpha_{t} - \alpha_{t}\alpha_{t-1}} + 1 - \alpha_{t}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Squaring the terms} \\ &= \sqrt{\alpha_{t}\alpha_{t-1}}\boldsymbol{x}_{t-2} + \sqrt{1 - \alpha_{t}\alpha_{t-1}}\boldsymbol{\epsilon}_{t-2} \qquad \text{Simplifying} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^{t}\alpha_{i}}\boldsymbol{x}_{0} + \sqrt{1 - \prod_{i=1}^{t}\alpha_{i}}\boldsymbol{\epsilon}_{0} \qquad \text{Let } \bar{\alpha}_{t} = \prod_{i=1}^{t}\alpha_{i} \\ &\sim \mathcal{N}(\boldsymbol{x}_{t}; \sqrt{\bar{\alpha}_{t}}\boldsymbol{x}_{0}, (1 - \bar{\alpha}_{t})\mathbf{I}) \qquad \mathbf{x}_{t} \text{ is now a Gaussian characterized by } \mathbf{x}_{0} \end{split}$$

Ranjay Krishna

Lecture 19 - 113

Takeaway from the previous slides:

 $oldsymbol{x}_t \sim \mathcal{N}(oldsymbol{x}_t; \sqrt{ar{lpha}_t}oldsymbol{x}_0, (1-ar{lpha}_t) \mathbf{I})$

We can instantly sample x_{t} given any input data x_{0}



Lecture 19 - 114



 $q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0)$



Lecture 19 - 115

$$q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) = rac{q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}, oldsymbol{x}_0)q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_0)}{q(oldsymbol{x}_t \mid oldsymbol{x}_0)}$$

Applying Bayes rule

Ranjay Krishna

Lecture 19 - 116

$$egin{aligned} q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) &= rac{q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}, oldsymbol{x}_0)}{q(oldsymbol{x}_t \mid oldsymbol{x}_0)} \ &= rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) oldsymbol{ extsf{I}})}{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) oldsymbol{ extsf{I}})} \end{aligned}$$

The first term is just a single forward diffusion process:

$$q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}) = \mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 117

$$egin{aligned} q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) &= rac{q(oldsymbol{x}_t \mid oldsymbol{x}_{t-1}, oldsymbol{x}_0) \overline{q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_0)}}{q(oldsymbol{x}_t \mid oldsymbol{x}_0)} &= rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-ar{lpha}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-ar{lpha}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-ar{lpha}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-ar{lpha}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-ar{lpha}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}} oldsymbol{x}_0, (1-oldsymbol{x}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{oldsymbol{x}_{t-1}} oldsymbol{x}_0, (1-oldsymbol{x}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I}) \overline{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{oldsymbol{x}_{t-1}} oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_t; \sqrt{lpha_t} oldsymbol{x}_{t-1}, (1-lpha_t) \mathbf{I})}{\mathbf{I}} oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}) \mathbf{I})} & = rac{\mathcal{N}(oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1}, (1-oldsymbol{x}_{t-1},$$

The second term is also a Gaussian using the formula we just derived:

$$oldsymbol{x}_t \sim \mathcal{N}(oldsymbol{x}_t; \sqrt{ar{lpha}_t}oldsymbol{x}_0, (1-ar{lpha}_t) \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 118

$$egin{aligned} q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{t}, oldsymbol{x}_{0}) &= rac{q(oldsymbol{x}_{t} \mid oldsymbol{x}_{0})q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{0})}{q(oldsymbol{x}_{t} \mid oldsymbol{x}_{0})} &= rac{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{lpha_{t}}oldsymbol{x}_{t-1}, (1-lpha_{t})\mathbf{I})\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{ar{lpha}_{t-1}}oldsymbol{x}_{0}, (1-ar{lpha}_{t-1})\mathbf{I}) & & \ \mathcal{N}(oldsymbol{x}_{t}; \sqrt{ar{lpha}_{t}}oldsymbol{x}_{0}, (1-ar{lpha}_{t})\mathbf{I}) & & \ \end{pmatrix} \end{aligned}$$

The third term is also a Gaussian using the same formula:

$$oldsymbol{x}_t \sim \mathcal{N}(oldsymbol{x}_t; \sqrt{ar{lpha}_t}oldsymbol{x}_0, (1-ar{lpha}_t) \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 119

$$egin{aligned} q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{t}, oldsymbol{x}_{0}) &= rac{q(oldsymbol{x}_{t} \mid oldsymbol{x}_{0})q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{0})}{q(oldsymbol{x}_{t} \mid oldsymbol{x}_{0})} &= rac{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{lpha_{t}}oldsymbol{x}_{t-1}, (1-lpha_{t})oldsymbol{I})\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{arlpha_{t-1}}oldsymbol{x}_{0}, (1-arlpha_{t})oldsymbol{I})}{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{arlpha_{t}}oldsymbol{x}_{0}, (1-arlpha_{t})oldsymbol{I})} &= rac{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{lpha_{t}}oldsymbol{x}_{t-1}, (1-lpha_{t})oldsymbol{I})\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{arlpha_{t-1}}oldsymbol{x}_{0}, (1-arlpha_{t})oldsymbol{I})}{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{arlpha_{t}}oldsymbol{x}_{0}, (1-arlpha_{t})oldsymbol{I})} &= rac{\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{arlpha_{t}}(1-arlpha_{t-1})oldsymbol{I})\mathcal{N}(oldsymbol{x}_{t-1}; \sqrt{arlpha_{t-1}}oldsymbol{I}) &= rac{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{arlpha_{t}}(1-oldsymbol{lpha}_{t-1})oldsymbol{I})}{\mathcal{N}(oldsymbol{x}_{t}; \sqrt{arlpha_{t}}oldsymbol{x}_{0}, (1-oldsymbol{lpha}_{t})oldsymbol{I}) &= rac{\mathcal{N}(oldsymbol{x}_{t-1}; oldsymbol{I})}{\sqrt{oldsymbol{lpha}_{t}} + \sqrt{oldsymbol{lpha}_{t-1}}(1-lpha_{t})oldsymbol{x}_{0})}, \underbrace{(1-lpha_{t})(1-oldsymbol{lpha}_{t-1})}{1-oldsymbol{lpha}_{t}}oldsymbol{I}) &= rac{\mathcal{N}(oldsymbol{x}_{t-1}; oldsymbol{I})}{\sqrt{oldsymbol{lpha}_{t}} + \sqrt{oldsymbol{lpha}}_{t-1}}(1-oldsymbol{lpha}_{t})oldsymbol{x}_{0})}, \underbrace{(1-lpha_{t})(1-oldsymbol{lpha}_{t-1})}{1-oldsymbol{lpha}_{t}} oldsymbol{I}) &= rac{\mathcal{N}(oldsymbol{x}_{t-1}; oldsymbol{A})}{\mathcal{N}(oldsymbol{x}_{t}, oldsymbol{x}_{0})}, \underbrace{(1-lpha_{t})(1-oldsymbol{lpha}_{t-1})}{1-oldsymbol{lpha}_{t}} oldsymbol{I}) &= rac{\mathcal{N}(oldsymbol{x}_{t-1}; oldsymbol{A})}{\mathcal{N}(oldsymbol{x}_{t}, oldsymbol{x}_{0})}, \underbrace{(1-oldsymbol{lpha}_{t})(1-oldsymbol{lpha}_{t-1})}{\mathbf{Y}(oldsymbol{A})} oldsymbol{A}) &= rac{\mathcal{N}(oldsymbol{x}_{t})}{\mathcal{N}(oldsymbol{x}_{t}, oldsymbol{A})}, \underbrace{(1-oldsymbol{A})}{\mathbf{Y}(oldsymbol{A})} &= rac{\mathcal{N}(olds$$

The product of these 3 Gaussian distributions simplify to a Gaussian as well!

Let's call its mean $\mu_q(\boldsymbol{x}_t, \boldsymbol{x}_0)$ and variance $\boldsymbol{\Sigma}_q(t)$

Ranjay Krishna

Lecture 19 - 120

Proof (out of scope for the class)

$$\begin{split} q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{t}, \boldsymbol{x}_{0}) &= \frac{q(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{t-1}, \boldsymbol{x}_{0})q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_{0})}{q(\boldsymbol{x}_{t} \mid \boldsymbol{x}_{0})} \\ &= \frac{\mathcal{N}(\boldsymbol{x}_{t}; \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1}, (1 - \alpha_{t})\mathbf{I})\mathcal{N}(\boldsymbol{x}_{t-1}; \sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{0}, (1 - \overline{\alpha}_{t-1})\mathbf{I})}{\mathcal{N}(\boldsymbol{x}_{t}; \sqrt{\overline{\alpha}_{t}}\boldsymbol{x}_{0}, (1 - \overline{\alpha}_{t})\mathbf{I})} \\ &\propto \exp\left\{-\left[\frac{(\boldsymbol{x}_{t} - \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1})^{2}}{2(1 - \alpha_{t})} + \frac{(\boldsymbol{x}_{t-1} - \sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{0})^{2}}{2(1 - \overline{\alpha}_{t-1})} - \frac{(\boldsymbol{x}_{t} - \sqrt{\overline{\alpha}_{t}}\boldsymbol{x}_{0})^{2}}{2(1 - \overline{\alpha}_{t})}\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\frac{(\boldsymbol{x}_{t} - \sqrt{\alpha_{t}}\boldsymbol{x}_{t-1})^{2}}{1 - \alpha_{t}} + \frac{(\boldsymbol{x}_{t-1} - \sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{0})^{2}}{1 - \overline{\alpha}_{t-1}} - \frac{(\boldsymbol{x}_{t} - \sqrt{\overline{\alpha}_{t}}\boldsymbol{x}_{0})^{2}}{1 - \overline{\alpha}_{t}}\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\frac{(-2\sqrt{\alpha_{t}}\boldsymbol{x}_{t}\boldsymbol{x}_{t-1} + \alpha_{t}\boldsymbol{x}_{t-1}^{2})}{1 - \alpha_{t}} + \frac{(\boldsymbol{x}_{t-1}^{2} - 2\sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{t-1}\boldsymbol{x}_{0})}{1 - \overline{\alpha}_{t-1}} + C(\boldsymbol{x}_{t}, \boldsymbol{x}_{0})\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[-\frac{2\sqrt{\alpha_{t}}\boldsymbol{x}_{t}\boldsymbol{x}_{t-1}}{1 - \alpha_{t}} + \frac{\alpha_{t}\boldsymbol{x}_{t-1}^{2}}{1 - \alpha_{t}} + \frac{2\sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{t-1}\boldsymbol{x}_{0}}{1 - \overline{\alpha}_{t-1}}\right)\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[(\frac{\alpha_{t}}{1 - \alpha_{t}} + \frac{1}{1 - \overline{\alpha}_{t-1}})\boldsymbol{x}_{t-1}^{2} - 2\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \overline{\alpha}_{t-1}}\right)\boldsymbol{x}_{t-1}\right]\right\} \\ &= \exp\left\{-\frac{1}{2}\left[\frac{\alpha_{t}(1 - \overline{\alpha}_{t-1}) + 1 - \alpha_{t}}{(1 - \alpha_{t})(1 - \overline{\alpha}_{t-1})}\boldsymbol{x}_{t-1}^{2} - 2\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\overline{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \overline{\alpha}_{t-1}}\right)\boldsymbol{x}_{t-1}\right]\right\} \end{aligned}$$

2025

Ra

Proof (out of scope for the class)

Rar

$$= \exp\left\{-\frac{1}{2}\left[\frac{\alpha_{t} - \bar{\alpha}_{t} + 1 - \alpha_{t}}{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})} \boldsymbol{x}_{t-1}^{2} - 2\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\bar{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t-1}}\right)\boldsymbol{x}_{t-1}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left[\frac{1 - \bar{\alpha}_{t}}{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})} \boldsymbol{x}_{t-1}^{2} - 2\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\bar{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t-1}}\right)\boldsymbol{x}_{t-1}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{1 - \bar{\alpha}_{t}}{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}\right)\left[\boldsymbol{x}_{t-1}^{2} - 2\frac{\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\bar{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t-1}}\right)}{\frac{1 - \bar{\alpha}_{t}}{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}}}\boldsymbol{x}_{t-1}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{1 - \bar{\alpha}_{t}}{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}\right)\left[\boldsymbol{x}_{t-1}^{2} - 2\frac{\left(\frac{\sqrt{\alpha_{t}}\boldsymbol{x}_{t}}{1 - \alpha_{t}} + \frac{\sqrt{\bar{\alpha}_{t-1}}\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t-1}}\right)(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t}}}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{1}{\frac{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t}}}\right)\left[\boldsymbol{x}_{t-1}^{2} - 2\frac{\sqrt{\alpha_{t}}(1 - \bar{\alpha}_{t-1})\boldsymbol{x}_{t} + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_{t})\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t}}}\boldsymbol{x}_{t-1}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{1}{\frac{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t}}}\right)\left[\boldsymbol{x}_{t-1}^{2} - 2\frac{\sqrt{\alpha_{t}}(1 - \bar{\alpha}_{t-1})\boldsymbol{x}_{t} + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_{t})\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t}}}\boldsymbol{x}_{t-1}\right]\right\}$$

$$= \exp\left\{-\frac{1}{2}\left(\frac{1}{\frac{(1 - \alpha_{t})(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_{t}}}\right)\left[\boldsymbol{x}_{t-1}^{2} - 2\frac{\sqrt{\alpha_{t}}(1 - \bar{\alpha}_{t-1})\boldsymbol{x}_{t} + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_{t})\boldsymbol{x}_{0}}{1 - \bar{\alpha}_{t}}}\boldsymbol{x}_{t-1}\right]\right\}$$

Let's go back to the Markovian VAE



We are ready to set up a simple intuitive loss function to train the decoder! Given an image x_0 :

We want to generate $p_{\theta}(\bm{x}_{t-1} \mid \bm{x}_t)$ to match the Gaussian we just derived: $q(\bm{x}_{t-1} \mid \bm{x}_t, \bm{x}_0)$

Ranjay Krishna

Lecture 19 - 123

You can show mathematically that:

 $q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0)$ Is also a Gaussian:

$$\mathcal{N}(oldsymbol{x}_{t-1}; \underbrace{rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t + \sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-ar lpha_t}, \underbrace{rac{(1-lpha_t)(1-ar lpha_{t-1})}{1-ar lpha_t}}_{oldsymbol{\Sigma}_q(t)} \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 124

You can show mathematically that:

 $q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0)$ Is also a Gaussian:

$$\mathcal{N}(oldsymbol{x}_{t-1}; \underbrace{rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t+\sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-ar lpha_t}, \underbrace{rac{(1-lpha_t)(1-ar lpha_{t-1})}{1-ar lpha_t}}_{oldsymbol{\Sigma}_q(t)} \mathbf{I})$$

And that:

$$oldsymbol{\mu}_q(oldsymbol{x}_t,oldsymbol{x}_0) = rac{\sqrt{lpha_t}(1-arlpha_{t-1})oldsymbol{x}_t+\sqrt{arlpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-arlpha_t} \ = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\epsilon}_0$$

Ranjay Krishna

Lecture 19 - 125

The loss function tries to match distributions



The loss function $\underset{\boldsymbol{\theta}}{\arg\min} \mathcal{D}_{\mathrm{KL}}(q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \mid\mid p_{\boldsymbol{\theta}}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t))$

Ranjay Krishna

Lecture 19 - 126

We can model $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t)$ as a Gaussian $p(x_0|x_1)$ $p(x_{t-1}|x_t)$ $p(x_t|x_{t+1})$ $p(x_{T-1}|x_T)$ x_0 x_{t-1} x_t x_{t+1} x_T . . . $q(x_1|x_0)$ $q(x_t|x_{t-1})$ $q(x_{t+1}|x_t)$ $q(x_T | x_{T-1})$ The loss function $rg\min \mathcal{D}_{ ext{KL}}(q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t, oldsymbol{x}_0) \mid\mid p_{oldsymbol{ heta}}(oldsymbol{x}_{t-1} \mid oldsymbol{x}_t))$

$$= \argmin_{\boldsymbol{\theta}} \mathcal{D}_{\mathrm{KL}}\left(\mathcal{N}\left(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_{q}, \boldsymbol{\Sigma}_{q}\left(t\right)\right) \mid\mid \mathcal{N}\left(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_{q}\left(t\right)\right)\right)$$

Ranjay Krishna

θ

Lecture 19 - 127

We can model $p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t)$ as a Gaussian $p(x_{t-1}|x_t)$ $p(x_t|x_{t+1})$ $p(x_0|x_1)$ $p(x_{T-1}|x_T)$ x_0 x_{t-1} x_t x_{t+1} x_T . . . $q(x_1|x_0)$ $q(x_t|x_{t-1})$ $q(x_{t+1}|x_t)$ $q(x_T | x_{T-1})$ The loss function N 11 1 1 Т 1 11 Т

$$egin{aligned} &rg\min_{oldsymbol{ heta}} \mathcal{D}_{ ext{KL}}(q(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{t}, oldsymbol{x}_{0}) \mid\mid p_{oldsymbol{ heta}}(oldsymbol{x}_{t-1} \mid oldsymbol{x}_{t})) \ &= rg\min_{oldsymbol{ heta}} rac{1}{2\sigma_{q}^{2}(t)} \Big[\|oldsymbol{\mu}_{oldsymbol{ heta}} - oldsymbol{\mu}_{q}\|_{2}^{2} \Big] \end{aligned}$$

Ranjay Krishna

Lecture 19 - 128

$$\begin{split} & \operatorname{Proof}\left(\operatorname{out} \operatorname{of} \operatorname{scope} \operatorname{for} \operatorname{class}\right) \\ & \underset{\theta}{\operatorname{arg\,min}} \mathcal{D}_{\operatorname{KL}}(q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \mid\mid p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t)) \\ & = \underset{\theta}{\operatorname{arg\,min}} \mathcal{D}_{\operatorname{KL}}(\mathcal{N}(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \mid\mid \mathcal{N}(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_{\theta}, \boldsymbol{\Sigma}_q(t))) \\ & = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2} \left[\log \frac{|\boldsymbol{\Sigma}_q(t)|}{|\boldsymbol{\Sigma}_q(t)|} - d + \operatorname{tr}(\boldsymbol{\Sigma}_q(t)^{-1}\boldsymbol{\Sigma}_q(t)) + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1}(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ & = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2} \left[\log 1 - d + d + (\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1}(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ & = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2} \left[(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T \boldsymbol{\Sigma}_q(t)^{-1}(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \\ & = \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2} \left[(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q)^T (\boldsymbol{\sigma}_q^2(t) \mathbf{I})^{-1}(\boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q) \right] \end{split}$$

Ranjay Krishna

Lecture 19 - 129

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_a^2(t)} \left[\| \boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_{q} \|_2^2 \right]$



Lecture 19 - 130

The loss we want to minimize is $\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\| \boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q \|_2^2 \right]$

From the previous slide, we got the mean from this:

$$\mathcal{N}(oldsymbol{x}_{t-1}; \underbrace{rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t+\sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-ar lpha_t}, \underbrace{rac{(1-lpha_t)(1-ar lpha_{t-1})}{1-ar lpha_t}}_{oldsymbol{\Sigma}_q(t)} \mathbf{I})$$

Ranjay Krishna

Lecture 19 - 131

The loss we want to minimize is
$$\arg \min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\| \boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q \|_2^2 \right]$$

From the previous slide, we got the mean from this:

$$\mathcal{N}(oldsymbol{x}_{t-1}; \underbrace{rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t+\sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-ar lpha_t}}_{\mu_q(oldsymbol{x}_t,oldsymbol{x}_0)}, \underbrace{rac{(1-lpha_t)(1-ar lpha_{t-1})}{1-ar lpha_t}}_{\mathbf{\Sigma}_q(t)} \mathbf{I})$$

So, we can write the mean to be: $\boldsymbol{\mu}_q(\boldsymbol{x}_t, \boldsymbol{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})\boldsymbol{x}_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)\boldsymbol{x}_0}{1 - \bar{\alpha}_t}$

Ranjay Krishna

Lecture 19 - 132

The loss we want to minimize is
$$\arg\min_{\theta} \frac{1}{2\sigma_q^2(t)} \left[\| \boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q \|_2^2 \right]$$

From the previous slide, we got the mean from this:

$$\mathcal{N}(oldsymbol{x}_{t-1}; \underbrace{rac{\sqrt{lpha_t}(1-ar lpha_{t-1})oldsymbol{x}_t+\sqrt{ar lpha_{t-1}}(1-lpha_t)oldsymbol{x}_0}{1-ar lpha_t}, \underbrace{rac{(1-lpha_t)(1-ar lpha_{t-1})}{1-ar lpha_t}}_{oldsymbol{\Sigma}_q(t)} \mathbf{I})$$

So, we can write the mean to be:

$$oldsymbol{\mu}_q(oldsymbol{x}_t,oldsymbol{x}_0) = rac{\sqrt{lpha_t}(1-arlpha_{t-1})oldsymbol{x}_t+\sqrt{arlpha}_{t-1}(1-lpha_t)oldsymbol{x}_0}{1-arlpha_t} \ = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\epsilon}_0$$

Ranjay Krishna

Lecture 19 - 133

Our neural network can predict noise instead!

$$oldsymbol{\mu}_q(oldsymbol{x}_t,oldsymbol{x}_0) \;= rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\epsilon}_0$$

We can also set our predicted mean to be:

$$oldsymbol{\mu}_{oldsymbol{ heta}}(oldsymbol{x}_t,t) = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\hat{\epsilon}}_{oldsymbol{ heta}}(oldsymbol{x}_t,t)$$

Why is this helpful?

Ranjay Krishna

Lecture 19 - 134

Our neural network can predict noise instead!

$$oldsymbol{\mu}_q(oldsymbol{x}_t,oldsymbol{x}_0) \;= rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\epsilon}_0$$

We can also set our predicted mean to be:

$$oldsymbol{\mu}_{oldsymbol{ heta}}(oldsymbol{x}_t,t) = rac{1}{\sqrt{lpha_t}}oldsymbol{x}_t - rac{1-lpha_t}{\sqrt{1-arlpha_t}}oldsymbol{\hat{\epsilon}}_{oldsymbol{ heta}}(oldsymbol{x}_t,t)$$

Why is this helpful? Because now our model needs to predict the noise that was injected, which turns out to be empirically more stable of an objective than predicting the image mean.

Ranjay Krishna

Lecture 19 - 135

The two loss objectives are equivalent

The loss function

$$\underset{\boldsymbol{\theta}}{\operatorname{arg\,min}} \mathcal{D}_{\mathrm{KL}}(q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \mid\mid p_{\boldsymbol{\theta}}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t))$$

$$= \underset{\boldsymbol{\theta}}{\operatorname{arg\,min}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \boldsymbol{\mu}_{\boldsymbol{\theta}} - \boldsymbol{\mu}_q \right\|_2^2 \right]$$
Instead of predicting the mean image values

Ranjay Krishna

Lecture 19 - 136

The two loss objectives are equivalent

$$\begin{aligned} & \underset{\theta}{\operatorname{arg\,min}} \mathcal{D}_{\operatorname{KL}}(q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \mid \mid p_{\theta}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t)) \\ &= \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2\sigma_q^2(t)} \boxed{\left\| \boldsymbol{\mu}_{\theta} - \boldsymbol{\mu}_q \right\|_2^2} & \text{Instead of predicting the mean image values} \\ &= \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2\sigma_q^2(t)} \frac{(1 - \alpha_t)^2}{(1 - \bar{\alpha}_t)\alpha_t} \boxed{\left\| \boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\theta}(\boldsymbol{x}_t, t) \right\|_2^2} & \text{The neural network can predict the added noise} \end{aligned}$$

Ranjay Krishna

Lecture 19 - 137

Proof: (out of scope)

$$\begin{split} & \arg\min_{\boldsymbol{\theta}} \mathcal{D}_{\mathrm{KL}}(q(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t, \boldsymbol{x}_0) \mid \mid p_{\boldsymbol{\theta}}(\boldsymbol{x}_{t-1} \mid \boldsymbol{x}_t)) \\ &= \arg\min_{\boldsymbol{\theta}} \mathcal{D}_{\mathrm{KL}}(\mathcal{N}(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_q, \boldsymbol{\Sigma}_q(t)) \mid \mid \mathcal{N}(\boldsymbol{x}_{t-1}; \boldsymbol{\mu}_{\boldsymbol{\theta}}, \boldsymbol{\Sigma}_q(t))) \\ &= \arg\min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1}{\sqrt{\alpha_t}} \boldsymbol{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t) - \frac{1}{\sqrt{\alpha_t}} \boldsymbol{x}_t + \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 \right\|_2^2 \right] \\ &= \arg\min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \boldsymbol{\epsilon}_0 - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \hat{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t) \right\|_2^2 \right] \\ &= \arg\min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} (\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t)) \right\|_2^2 \right] \\ &= \arg\min_{\boldsymbol{\theta}} \frac{1}{2\sigma_q^2(t)} \left[\left\| \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} (\boldsymbol{\epsilon}_0 - \hat{\boldsymbol{\epsilon}}_{\boldsymbol{\theta}}(\boldsymbol{x}_t, t)) \right\|_2^2 \right] \end{split}$$

Ranjay Krishna

Lecture 19 - 138

The denoising architecture



Time representation: sinusoidal positional embeddings.

Ranjay Krishna

Lecture 19 - 139

How do we sample a new image?

Sample $x_T \sim \mathcal{N}(0, I)$ For $t = T \dots 1$ do Predict $\hat{\epsilon}_t = p_{\theta}(x_t)$ $\mu_{t-1} = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}\sqrt{\alpha_t}} \sqrt{1 - \bar{\alpha}_{t-1}} \hat{\epsilon}_t$ Sample $x_{t-1} \sim \mathcal{N}(\mu_{t-1}, \Sigma_{t-1})$ Return x_0



Reverse denoising process (generative)

Ranjay Krishna

Lecture 19 - 140

How is the time step inputted:

Ranjay Krishna

Time representation: sinusoidal positional embeddings.

Added in using: $AdaGN(h, y) = y_s \operatorname{GroupNorm}(h) + y_b$

- *h* is the intermediate activations of the residual block following the first convolution in each layer,

Lecture 19 - 141

- $y = [y_{s'}, y_{b}]$ is obtained from a linear projection of the timestep



Text-conditioned generation



Ranjay Krishna

Lecture 19 - 142

Application: panorama generation



Ranjay Krishna

Lecture 19 - 143

Application: super-resolution

Learn a superresolution diffusion model conditioned on a low resolution image. *y* is a low resolution input image, x is a high resolution output image

 $\mathbb{E}_{\mathbf{x},\mathbf{y}} \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0},\mathbf{I})} \mathbb{E}_t ||\epsilon_{\theta}(\mathbf{x}_t,t;\mathbf{y}) - \epsilon||_p^p$



Saharia et al., Image Super-Resolution via Iterative Refinement, 2021

Ranjay Krishna

Lecture 19 - 144
Application: super resolution

Natural Image Super-Resolution $64 \times 64 \rightarrow 256 \times 256$



Saharia et al., Image Super-Resolution via Iterative Refinement, 2021

Ranjay Krishna

Lecture 19 - 145

Application: image editing



Meng et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022

Ranjay Krishna

Lecture 19 - 146

Latent diffusion models: perform diffusion over latent VAE encodings



Rombach et al. High-Resolution Image Synthesis with Latent Diffusion Models ArXiv 2022

Ranjay Krishna

Lecture 19 - 147

Stable diffusion - from Stability Al

- Open sourced diffusion model main model used for research
- Produces 512x512 images
- UNet with 860M params
- ViT-L text encoder with 123M params
- Fits in 10GB VRAM fits on most GPUs



Rombach et al. High-Resolution Image Synthesis with Latent Diffusion Models ArXiv 2022

Ranjay Krishna

Lecture 19 - 148

Imagen - Google

Combines:

- Latent diffusion model
- text conditioning
- 2 super-resolution models

To produce high quality 1024x1024 images

Saharia et al., "Photorealistic Text-to-Image Diffusion Models with Deep Language Understanding", arXiv 2022.

Ranjay Krishna



"A Golden Retriever dog wearing a blue checkered beret and red dotted turtleneck."



Imagen examples



A dragon fruit wearing karate belt in the snow.

A relaxed garlic with a blindfold reading a newspaper while floating in a pool of tomato soup.

A photo of a Shiba Inu dog with a backpack riding a bike. It is wearing sunglasses and a beach hat.

Ranjay Krishna

Lecture 19 - 150

Last year: Sora video diffusion model

https://openai.com/sora

How did they do it?

- More data (unknown data source)
- Replaced U-Net architecture with transformers

Ranjay Krishna

Lecture 19 - 151

Q. Which ones are VAEs good at?

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

Ranjay Krishna

Lecture 19 - 152



VAEs are bad at generating high quality samples

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

Ranjay Krishna

Lecture 19 - 153

Q. Which ones are GANs good at?

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

Ranjay Krishna

Lecture 19 - 154

GANs suffer from mode collapse

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

Ranjay Krishna

Lecture 19 - 155

Q. Which ones are Diffusion models good at?

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

Ranjay Krishna

Lecture 19 - 156

Diffusion models are bad at sampling fast.

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

Ranjay Krishna

Lecture 19 - 157

End of course

What should you do next? Pursue a career in deep learning Start deep learning research



Lecture 19 - 158

