

Deep Learning

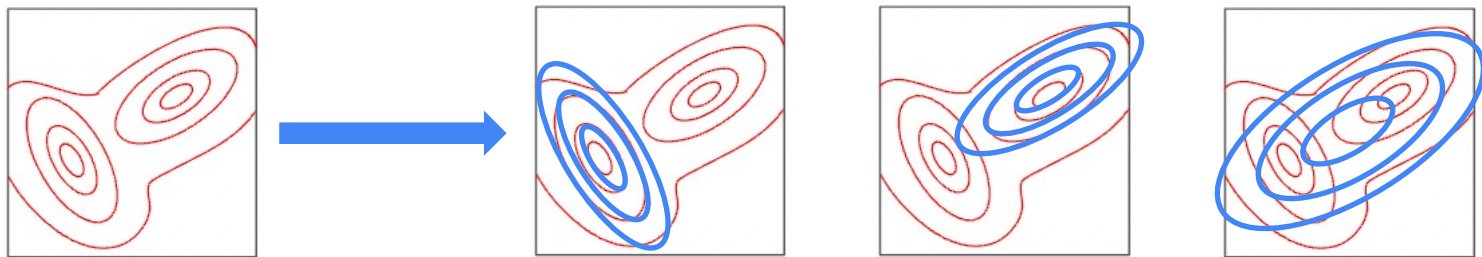
Diffusion Models and Normalizing Flows

11-785 - Fall 2023

Abuzar Khan

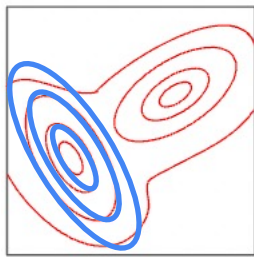
Open Question

1. Say we have a data distribution \mathbf{p} that is a mixture of two 2D gaussians as shown below in red. We want to approximate this with one gaussian estimate \mathbf{q} using KL-divergence. Which of the following three will result from optimizing $\mathbf{D}_{\text{KL}}(p || q)$
2. and which from $\mathbf{D}_{\text{KL}}(q || p)$?

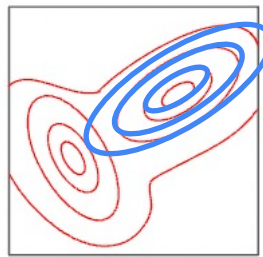


Open Question

1. Say we have a data distribution **p** that is a mixture of two 2D gaussians as shown below in red. We want to approximate this with one gaussian estimate **q** using KL-divergence. Which of the following three will result from optimizing $\mathbf{D}_{\text{KL}}(p || q)$
2. and which from $\mathbf{D}_{\text{KL}}(q || p)$?



$$\mathbf{D}_{\text{KL}}(q || p)$$



$$\mathbf{D}_{\text{KL}}(q || p)$$



$$\mathbf{D}_{\text{KL}}(p || q)$$

Background

1. Generative Models and Discriminative models
2. Autoencoders
3. Variational Autoencoders
 1. Reparameterization trick
 2. ELBO

Sandcastles

How to create a sandcastle:

Step 1: Take a sandcastle

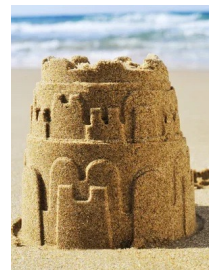
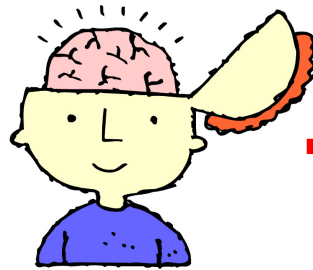
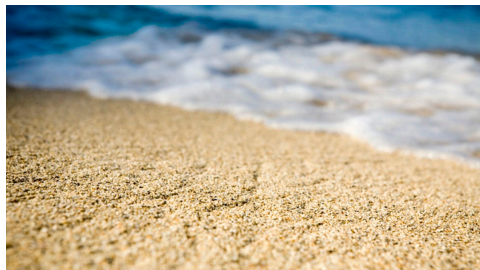
Step 2: Destroy the sandcastle

Step 3: Remember how you destroyed the sandcastle

Step 4: Reverse the process

Key Idea

Once you know how to reconstruct sandcastles, you can start with some different “sand”, apply this process, and end up with a different “sandcastle”



Part 1

Diffusion Models

ELBO Recap

Why use ELBO?

Directly maximizing $p(x)$ is very difficult:

- it involves either marginalizing over the entire latent space \mathbf{Z} (intractable for complex models) OR
- It involves having access to the ground truth latent encoder $p(z|x)$

ELBO:

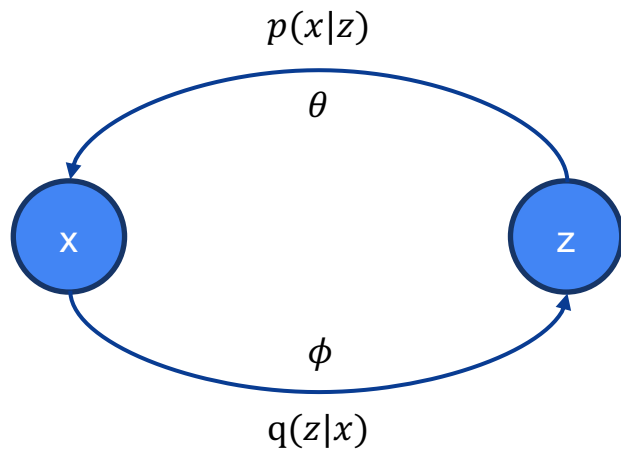
$$\log(p(x)) \geq \mathbb{E}_{q_{\phi}(z|x)} \left[\log \frac{p(x, z)}{q_{\phi}(z|x)} \right]$$

Question: Why does the \geq show up here? \rightarrow With the derivation in the appendix, we see a $D_{KL}(q_{\phi}(z|x) || p(z|x))$ term show up which is always ≥ 0 .

Applying chain-rule of probabilities:

$$ELBO = \underbrace{\mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)]}_{\text{Reconstruction}} - \underbrace{D_{KL}(q_{\phi}(z|x) || p(z))}_{\text{Prior matching}}$$

Variational Autoencoder Recap



Latent variable sampling: $z \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))$

Reparameterization trick: $z = \mu_\phi(x) + \sigma_\phi(x) \odot \epsilon, \epsilon \sim \mathcal{N}(0, I)$

Training:

- Jointly optimize θ and ϕ
- Maximize **ELBO**

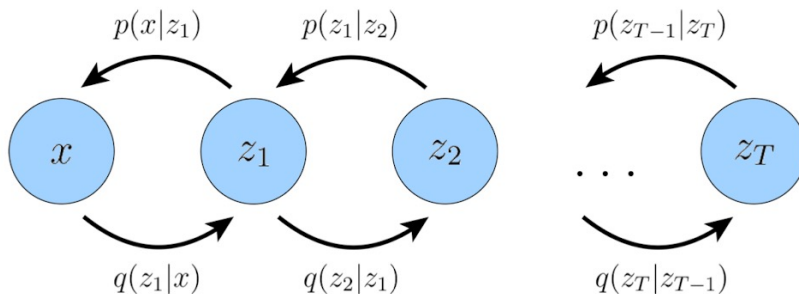
Empirically, we found that two things make VAEs work really well:

1. Increasing the depth of the networks
2. Introducing a hierarchy of latent variables (latent variables of latent variables)

$x \leftarrow z_1 \leftarrow z_2 \leftarrow \dots \leftarrow z_T$, such that each latent is conditioned on all previous latents.

We are particularly interested in such HAVEs that where the process is a **Markovian chain - MHVAE**

Markovian Hierarchical Variational Autoencoder



Joint probability: $p(x, z_{1:T}) = p(z_T) p_\theta(x | z_1) \prod_{t=2}^T p_\theta(z_{t-1} | z_t)$

Posterior probability: $q_\phi(z_{1:T} | x) = q_\phi(z_1 | x) \prod_{t=2}^T q_\phi(z_t | z_{t-1})$

Updated ELBO:

$$\log(p(x)) \geq \mathbb{E}_{q_\phi(z_{1:T} | x)} \left[\log \frac{p(x, z_{1:T})}{q_\phi(z_{1:T} | x)} \right]$$

Diffusion Models

Diffusion models are essentially **MHVAEs** with **3 restrictions**:

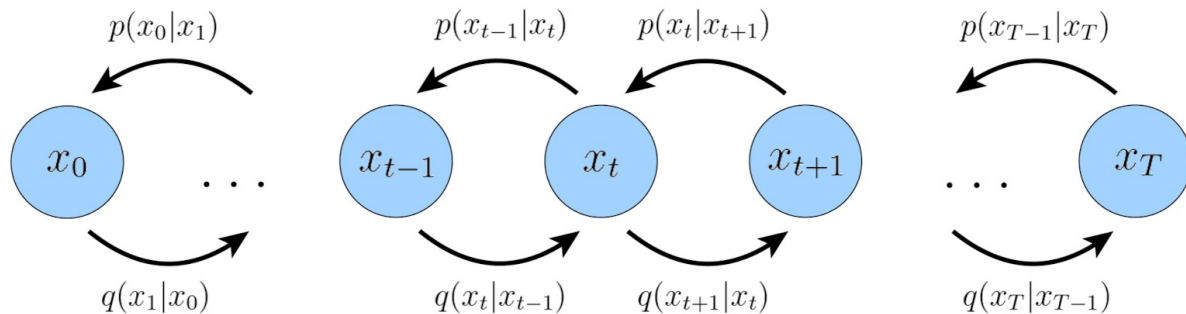
1. Latent dimension is the same as the data dimension
2. The encoder has no parameters to be learnt. It is defined to be a linear gaussian such that the t^{th} gaussian is centered around the previous latent z_{t-1}
3. The parameters for the gaussians are scheduled such that the final latent is a standard gaussian.

$$z_T \sim \mathcal{N}(z_T; 0, \mathbf{I})$$

The first restriction allows for some mild abuse of notation:

$$q_\phi(x_{1:T} | x_0) = \prod_{t=1}^T q_\phi(x_t | x_{t-1})$$
$$p(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

Diffusion Models



The first restriction allows for some mild abuse of notation:

$$q_\phi(x_{1:T} | x_0) = \prod_{t=1}^T q_\phi(x_t | x_{t-1})$$

$$p(x_{0:T}) = p(x_T) \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

Diffusion Models – Diffusion Process

Following the second restriction, we now define the linear gaussian for the encoding (diffusion) process:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \mu_t(x_{t-1}), \Sigma_t \mathbf{I})$$
$$\mu_t(x_{t-1}) = \sqrt{1 - \beta_t} x_{t-1}, \quad \Sigma_t = \beta_t$$

We additionally define $\alpha_t = 1 - \beta_t$.

β_t is defined to preserve variance across the diffusion steps.

We can now write

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

Using the reparameterization trick:

$$x_t = \sqrt{\alpha_t} x_{t-1} + (\sqrt{1 - \alpha_t}) \epsilon, \quad \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

$$= \sqrt{\alpha_t \alpha_{t-1}} x_{t-2} + (\sqrt{1 - \alpha_t \alpha_{t-2}}) \epsilon$$

$$= \sqrt{\alpha_t \alpha_{t-1} \alpha_{t-2}} x_{t-3} + (\sqrt{1 - \alpha_t \alpha_{t-2} \alpha_{t-3}}) \epsilon$$

$$= \sqrt{\alpha_t \alpha_{t-1} \dots \alpha_1} x_{(0)} + (\sqrt{1 - \alpha_t \alpha_{t-2} \dots \alpha_1}) \epsilon$$

$$= \sqrt{\bar{\alpha}_t} x_{(0)} + (\sqrt{1 - \bar{\alpha}_t}) \epsilon$$

This takes us from time step 0 to t
in **one step!**

From the third restriction, we get

$$\alpha_T \rightarrow 0$$

Sum of two gaussians is another gaussian with mean as the sum of the two means and variance as the sum of the two variances.

$$(1 - \alpha_t) \epsilon \rightarrow \mathcal{N}(\epsilon; 0, 1 - \alpha_t \mathbf{I})$$

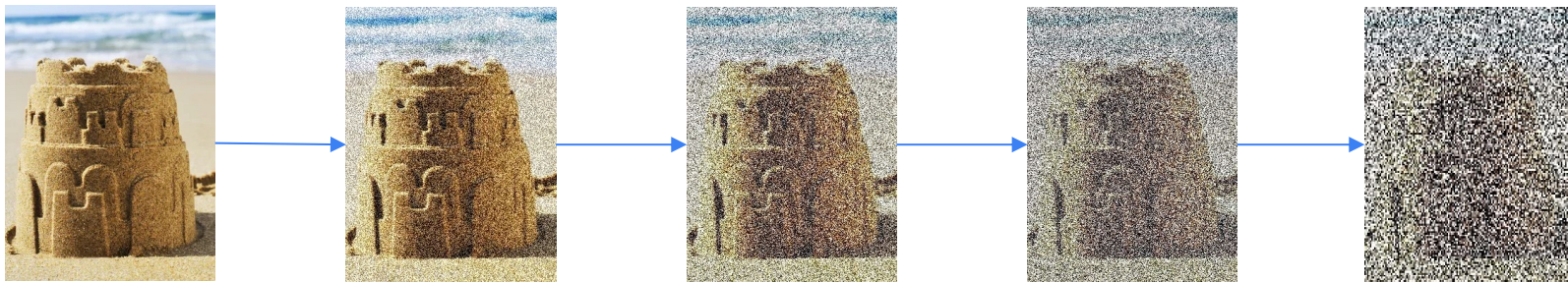
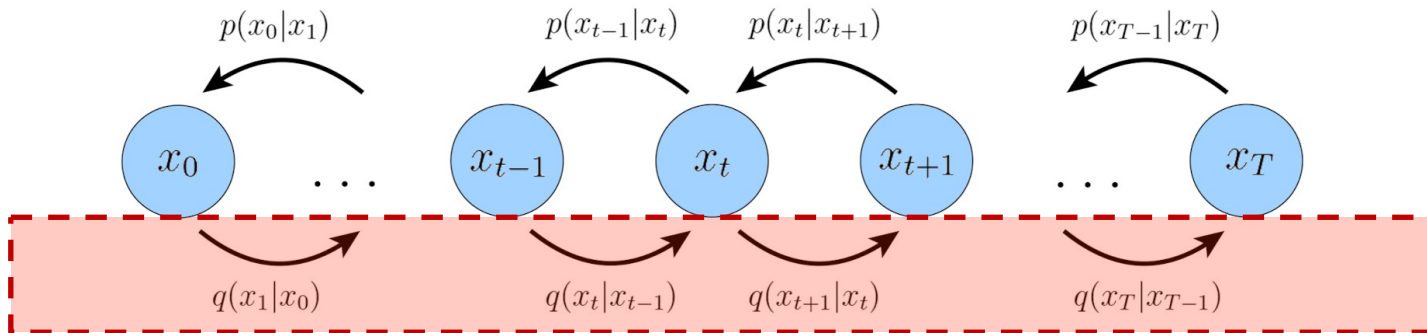
Define

$$\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$$

Diffusion Models – Diffusion Process

$$q(x_t|x_{t-1}) = \sqrt{\alpha_t}x_{(t-1)} + (\sqrt{1 - \alpha_t})\epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

This formulation essentially paints a picture of this process to be incrementally adding noise till we reach x_T which is defined to be pure noise.



Diffusion Models – Generative Process

From the third assumption, we can write the exact prior on the final step x_T :

$$p(x_T) = \mathcal{N}(x_T; 0, I)$$

For all other steps, we can write a learned distribution:

$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_t I)$$

Neural Network: **U-Net**
Denoising network

Exactly tractable
variance

$$p_{\theta}(x_{0:T}) = p(x_T) \prod_{t=1}^T p_{\theta}(x_{t-1}|x_t)$$

Diffusion Models – Updated ELBO

*Derivation in appendix!

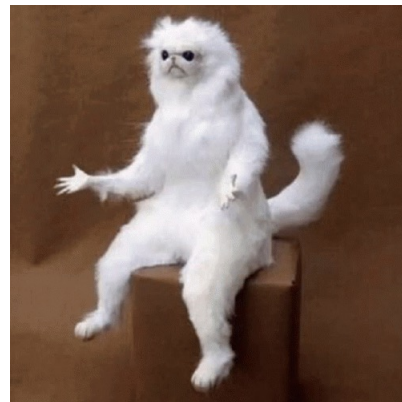
$$\begin{aligned}\log p(x) &= \log \int p(x_{0:T}) dx_{0:T} \\ &\quad \dots * \\ &= \mathbb{E}_{q(x_1 | x_0)} [\log p_\theta(x_0 | x_1)] - \mathbb{E}_{q(x_{T-1} | x_0)} [D_{KL}(q(x_T | x_{T-1}) || p(x_T))] \\ &\quad - \sum_{t=1}^{T-1} \mathbb{E}_{q(x_{t-1}, x_{t+1} | x_0)} [D_{KL}(q(x_t | x_{t-1}) || p_\theta(x_t | x_{t+1}))]\end{aligned}$$

↑ ↑

This has **2** random variables for each **t**, this makes the computation slightly hard. We would prefer for there to be need for just **1**!

We can arbitrarily modify the diffusion process distribution to

$$q(x_t | x_{t-1}, x_0) = \frac{q(x_{t-1} | x_t, x_0) q(x_t | x_0)}{q(x_{t-1} | x_0)}$$



Diffusion Models – Updated ELBO

*Derivation in appendix!



$$\log p(x) = \log \int p(x_{0:T}) dx_{0:T}$$

... *

$$= \mathbb{E}_{q(x_1 | x_0)} [\log p_\theta(x_0 | x_1)] - D_{KL}(q(x_T | x_0) || p(x_T)) - \sum_{t=2}^T \mathbb{E}_{q(x_t | x_0)} [D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t))]$$

Reconstruction

Prior matching

Denoising

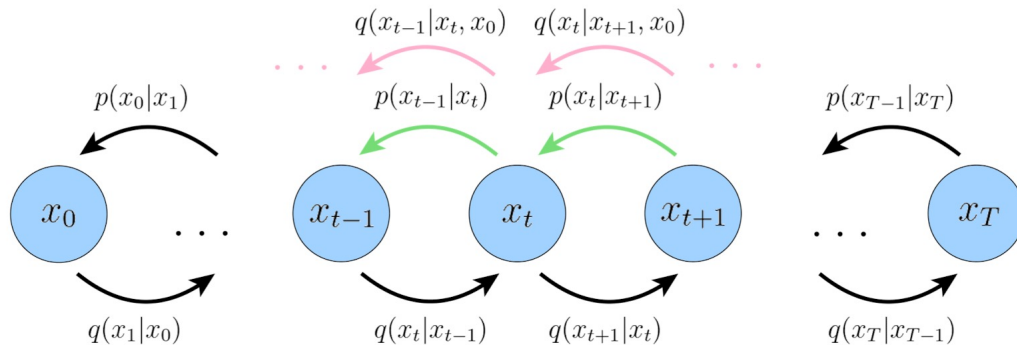
- **Reconstruction**: Reconstruction from least noisy version (**hyperparameter choice can make this arbitrarily small**)
- **Prior matching**: Moving the posterior closer to the true prior on the final noisy step (**0 for diffusion models**)
- **Denoising**: Divergence between approximate denoising (p_θ) and true denoising (q) steps

$q(x_{t-1} | x_t, x_0)$ is **tractable** and can be calculated **exactly** without any approximation:

$$q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1}; \bar{\mu}_t, \Sigma_t I)$$

$$\bar{\mu}_t = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})x_t + \sqrt{\bar{\alpha}_{t-1}}(1 - \alpha_t)x_0}{1 - \bar{\alpha}_t}, \quad \Sigma_t = \frac{(1 - \alpha_t)(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}$$

Diffusion Models – Loss formulation



Loss can focus on the denoising term. Decomposing for each timestep, we can have the t^{th} loss term:

$$L_t = D_{KL}(q(x_{t-1}|x_t, x_0) || p_{\theta}(x_{t-1}|x_t)) + \mathcal{C}$$

Since both inputs of the divergence are gaussians, this further simplifies to:

$$L_t = \mathbb{E}_q \left[\frac{1}{2\Sigma_t} ||\bar{\mu}_t - \mu_{\theta}(x_t, t)||^2 \right] + \mathcal{C}$$

Diffusion Models – Loss formulation

Further, we have $x_t = \sqrt{\alpha_t}x_{(t-1)} + (\sqrt{1 - \alpha_t})\epsilon$, $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ from definition

This lets us rewrite the true mean of the denoising process as:

$$\bar{\mu}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \bar{\alpha}_t}} \epsilon \right)$$

We can also write the predicted mean as:

$$\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{(1 - \alpha_t)}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(x_t, t) \right)$$

This lets us reformulate the loss to present a noise prediction problem:

$$L_{t-1} = \mathbb{E}_{x_0, \epsilon} \left[\frac{(1 - \alpha_t)^2}{2\Sigma_t \alpha_t (1 - \bar{\alpha}_t)} \|\epsilon - \epsilon_{\theta}(x_t, t)\|^2 \right] + C$$

Diffusion Models – Training and Inference

Algorithm 1 Training

```
1: repeat  
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$   
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$   
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
5:   Take gradient descent step on  $\epsilon$  on  $\mathbf{x}_t$   
         $\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2$   
6: until converged
```

Algorithm 2 Sampling

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

How do we tell the model what timestep we are on?

Temporal encodings in the form of sinusoids (or anything, really)

Diffusion models - Summary

- Diffusion models are **Markovian Hierarchical VAEs** with extra restrictions
- The loss is the vanilla VAE ELBO loss with an added denoising term
- The encoder has **0 parameters**
- The true denoising posterior can be **exactly calculated**
- The problem can be reformulated as a noise prediction problem
- There's a ton of math underlying a rather simple intuition

Part 2

Normalizing Flows

Sandcastles

How to create a sandcastle:

Step 1: Take a sandcastle

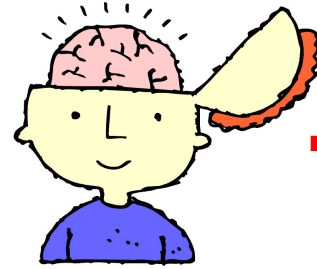
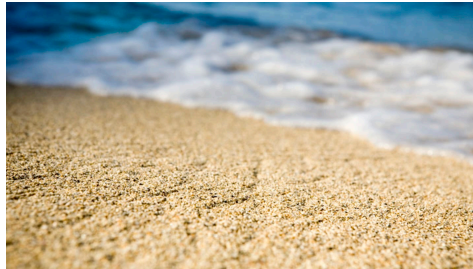
Step 2: Destroy the sandcastle

Step 3: Remember how you destroyed the sandcastle

Step 4: Reverse the process

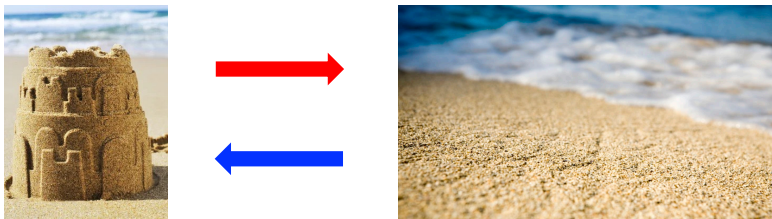
Key Idea

Once you know how to reconstruct sandcastles, you can start with some different “sand”, apply this process, and end up with a different “sandcastle”

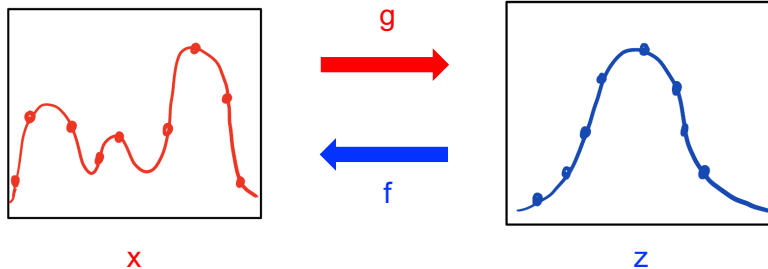


Normalizing Flows – Motivation

- In VAEs we are faced with an intractable likelihood calculation
- We use an ELBO instead as a surrogate objective to MLE
- What if we wanted to do MLE exactly?
- That would require us to go from sandcastle to sand, and back, without any approximation or estimation!

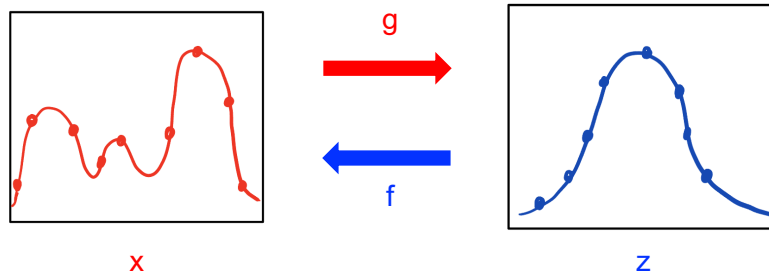


We would need for all the steps we do to be **invertible!**



It follows that
 $f^{-1} = g$

Normalizing Flows – Log likelihood



Bijection (and invertibility) allow us to directly compute the likelihood:

$$\int p_x(x)dx = \int p_z(g(x))dz$$

In multiple dimensions,
we generalize to the
determinant of the
Jacobian

$$p_x(x) = p_z(g(x)) \left| \frac{dg(x)}{dx} \right| \rightarrow p_z(g(x)) |det. J(g(x))|$$

$$\log p_x(x) = \log p_z(g(x)) + \log |det. J(g(x))|$$

Intuitively

$z = g(x)$ determines
where a point in x-space
maps to z-space (where
to move grains of sand)

$|det. J(g(x))|$ describes
how much probability
mass (sand) gets moved
in a local neighborhood.

Normalizing Flows – Closer look at the Jacobian

$$\mathbf{z} = g(\mathbf{x}) = f^{-1}(\mathbf{x})$$

$$\mathbf{J}_{\mathbf{x}} g(\mathbf{x}) = \begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_K} \\ \vdots & & \vdots \\ \frac{\partial z_K}{\partial x_1} & \dots & \frac{\partial z_K}{\partial x_K} \end{bmatrix}$$

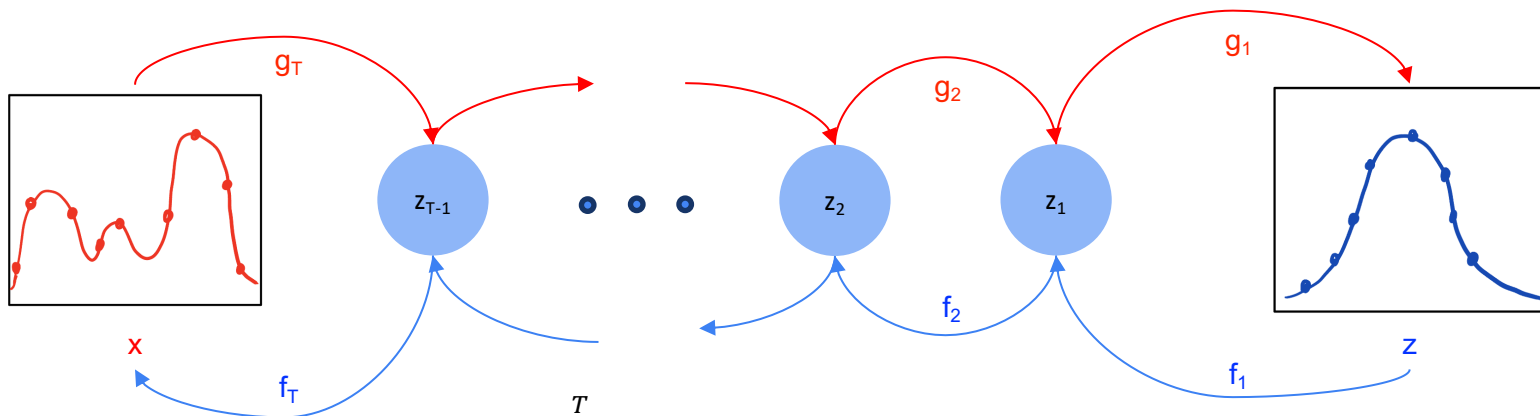
Normalizing Flows – Compositions

Bijections allow for composing several functions together!

This follows that we can now define:

$$z \sim p(z)$$

$$x = f_T \circ f_{T-1} \circ \dots \circ f_1(z)$$



$$\log p_x(x) = \log p_z(z_0) + \sum_{t=1}^T \log |\det J_{z_t}(g(z_{t-1}))|, \quad z_T = x, z_0 = z$$

Inverse: $z = g_1 \circ g_2 \circ \dots \circ g_T(x)$

Normalizing Flows – Characteristics

For a good (efficient) flow, we must have functions (steps) that are:

1. Expressive
2. Invertible
3. Offer cheap to compute Jacobian determinants

Computing a determinant is a **cubic** operation, but some special cases of matrices can make it very cheap.

Especially, **diagonal** matrices:

For a diagonal matrix, the determinant is simply the product of its diagonal elements. Same applies for any triangular matrix!

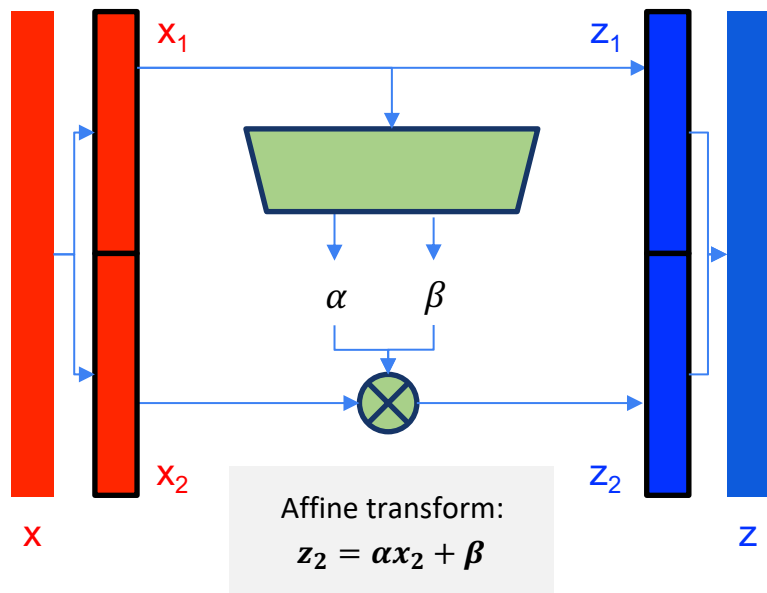
In [linear algebra](#), a **diagonal matrix** is a [matrix](#) (usually a [square matrix](#)) in which the entries outside the [main diagonal](#) (\searrow) are all zero. The diagonal entries themselves may or may not be zero. Thus, the matrix $D = (d_{ij})$ with n columns and n rows is diagonal if:

$$d_{i,j} = 0 \text{ if } i \neq j \quad \forall i, j \in \{1, 2, \dots, n\}$$

For example, the following matrix is diagonal:

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 4 & 0 \\ 0 & 0 & -2 \end{bmatrix}$$

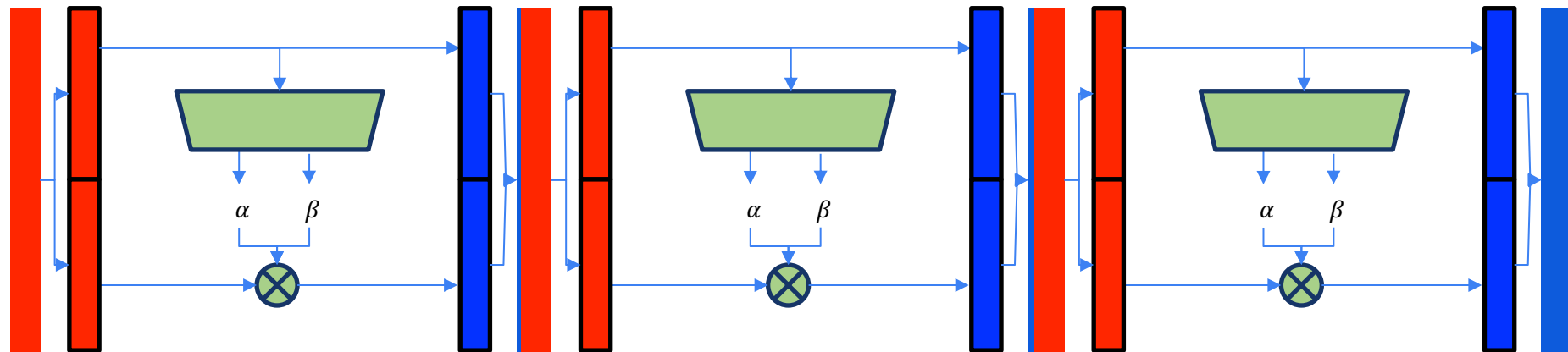
Normalizing Flows – Construction



$$J_z(z)$$

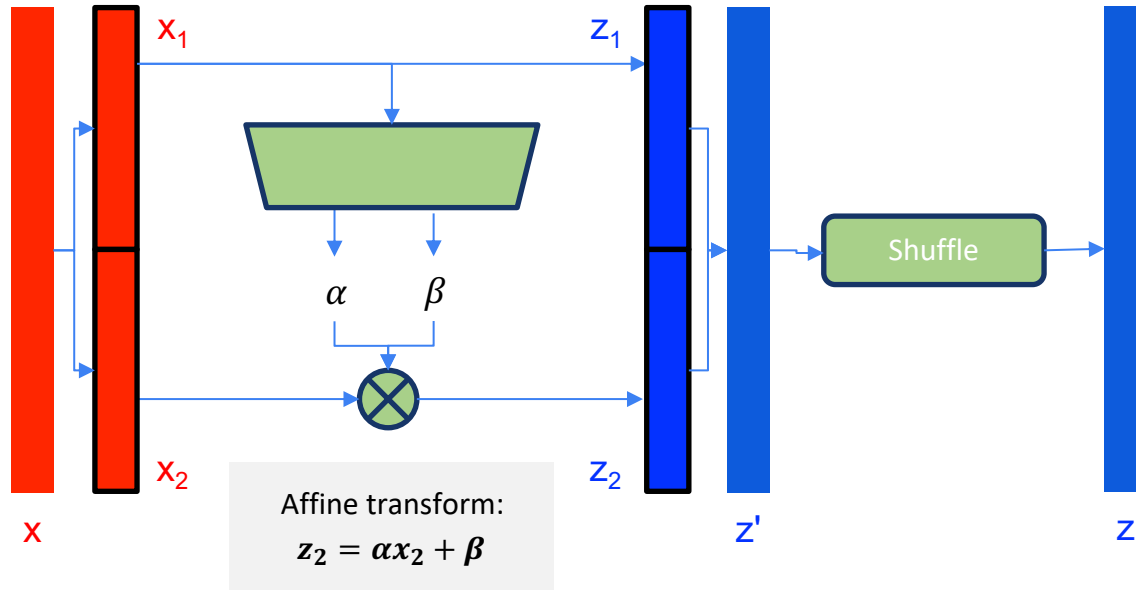
$$J_z(z) = \begin{bmatrix} x_1 & x_2 \\ z_1 & z_2 \end{bmatrix} = \begin{bmatrix} I & 0 \\ ? & \alpha I \end{bmatrix}$$

Normalizing Flows – Composition



Is there a problem here?

Normalizing Flows – Construction with a shuffle



This is the most popular type of flow called as **Coupling Flow** – Used in implementations such as NICE and GLOW

Normalizing Flows – In practice for images

- Multiscale architecture
- Split along channels
- Employ CNNs
- Perform permutations using 1x1 Conv layers

} GLOW!

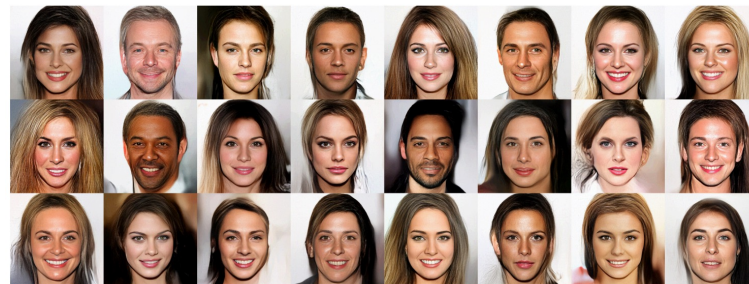


Figure 4: Random samples from the model, with temperature 0.7.



Figure 5: Linear interpolation in latent space between real images.

References

- Denoising Diffusion Probabilistic Models, <https://arxiv.org/abs/2006.11239>
- Diffusion Models: A Comprehensive Survey of Methods and Applications, <https://arxiv.org/abs/2209.00796>
- Understanding Diffusion Models: A Unified Perspective, <https://arxiv.org/abs/2208.11970>
- GLOW, <https://proceedings.neurips.cc/paper/2018/file/d139db6a236200b21cc7f752979132d0-Paper.pdf>
- Variational inference with normalizing flows, <https://arxiv.org/abs/1505.05770>
- Normalizing Flows: An Introduction and Review of Current Methods, <https://arxiv.org/abs/1908.09257>

- Didrik Nielsen's lecture, <https://www.youtube.com/watch?v=2tVHbcUP9b8>
- Hans van Gorp's lecture, <https://www.youtube.com/watch?v=yxVcnuRrKqQ&t=17s>
- Tim Salimans' lecture, <https://www.youtube.com/watch?v=pea3sH6orMc>

Appendix