

Deep Generative Models

Rianne van den Berg, principal researcher @Microsoft Research Amsterdam

Previously at Google Brain & University of Amsterdam



Outline

1. Variational auto-encoders
 - a. Variational inference
 - b. Variational auto-encoders
2. Normalizing flows
 - a. Normalizing flows for variational inference
 - b. Generative normalizing flows
3. Denoising diffusion models
4. Generative models for lossless compression
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Outline

- 1. Variational auto-encoders**
 - a. Variational inference
 - b. Variational auto-encoders
- 2. Normalizing flows**
 - a. Normalizing flows for variational inference
 - b. Generative normalizing flows
- 3. Denoising diffusion models**
- 4. Generative models for lossless compression**
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Approximate variational inference

Jordan et al, 1998; Jaakkola and Jordan 2000; Jaakkola 2001; Wainwright and Jordan 2008

Variational inference: a review for statisticians, Blei et al. 2018

“The goal of variational inference is to approximate a conditional density of latent variables given observed variables.”

Observed data: $X = \{x_1, \dots, x_N\}$

We assume that the data was produced by the generative model

$$p(z, x) = p(z)p(x|z)$$

Latent variables / model parameters

Approximate variational inference

Example: Bayesian mixture of K Gaussians

K cluster means: $\mu_k \sim p(\mu_k) = \mathcal{N}(0, \sigma^2) \quad k = 1, \dots, K$

N cluster assignments: $\mathbf{z}_i \sim p(\mathbf{z}_i) = \text{Categorical}\left(\frac{1}{K}, \dots, \frac{1}{K}\right) \quad i = 1, \dots, N$

N datapoints: $x_i \sim p(x_i | \mathbf{z}_i, \boldsymbol{\mu}) = \mathcal{N}(x_i | z_i^T \boldsymbol{\mu}, 1) \quad i = 1, \dots, N$

$$p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) = p(\boldsymbol{\mu}) \prod_{i=1}^N p(\mathbf{z}_i) p(x_i | \mathbf{z}_i, \boldsymbol{\mu})$$

Goal: infer $p(\boldsymbol{\mu}, \mathbf{Z} | \mathbf{X}) = p(\mathbf{X}, \mathbf{Z}, \boldsymbol{\mu}) / p(\mathbf{X})$

Evidence: intractable

Approximate variational inference

Pick a family of distributions Q , and find the best candidate $q(z) \in Q$ such that

$$q^*(z) = \arg \min_{q(z) \in Q} KL[q(z) || p(z|x)]$$
$$KL[q(z) || p(z|x)] = \mathbb{E}[\log q(z)] - \mathbb{E}[\log p(z|x)]$$

But this KL is not tractable because we don't know $p(z|x)$...

$$KL[q(z) || p(z|x)] = \mathbb{E}[\log q(z)] - \mathbb{E}[\log p(z, x)] + \log p(x)$$

Instead we can maximize an alternative objective:

$$ELBO(q) = \mathbb{E}[\log p(z, x)] - \mathbb{E}[\log q(z)]$$

Approximate variational inference

$$ELBO(q) = \mathbb{E}[\log p(z, x)] - \mathbb{E}[\log q(z)]$$

ELBO: Evidence lower bound

$$\log p(x) = KL[q(z) || p(z|x)] + ELBO(q) \rightarrow ELBO \leq \log p(x)$$

Optimizing the ELBO with respect to $q(z)$ ensures that

1. $q(z)$ approximates the unknown $p(z|x)$
2. $ELBO(q)$ becomes a tighter lower bound to $\log p(x)$

Mean-field variational inference

Mean-field variational family: all latent variables are independent and each variable has a distinct factor describing its distribution

$$q(z) = \prod_{j=1}^m q_j(z_j)$$

K cluster means:

$$\mu_k \sim p(\mu_k) = \mathcal{N}(0, \sigma^2) \quad k = 1, \dots, K$$

N cluster assignments:

$$\mathbf{z}_i \sim p(\mathbf{z}_i) = \text{Categorical}\left(\frac{1}{K}, \dots, \frac{1}{K}\right) \quad i = 1, \dots, N$$

N datapoints:

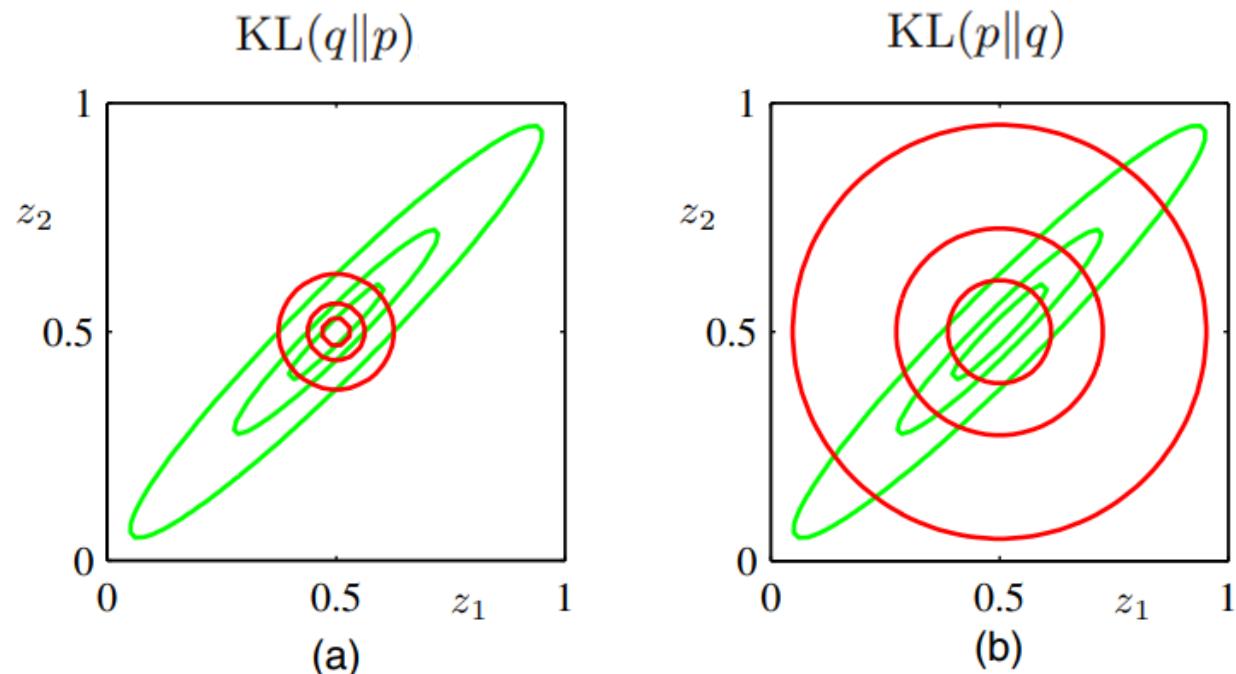
$$x_i \sim p(x_i | \mathbf{z}_i, \boldsymbol{\mu}) = \mathcal{N}(x_i | z_i^T \boldsymbol{\mu}, 1) \quad i = 1, \dots, N$$

$$p(\boldsymbol{\mu}, \mathbf{Z} | X) = ?$$

$$q(\boldsymbol{\mu}, \mathbf{Z}) = \prod_{k=1}^K q(\mu_k | m_k, s_k^2) \prod_{i=1}^N q(z_i | \phi_i)$$

KL and reverse KL: zero avoiding and zero forcing

Figure 10.2 Comparison of the two alternative forms for the Kullback-Leibler divergence. The green contours corresponding to 1, 2, and 3 standard deviations for a correlated Gaussian distribution $p(\mathbf{z})$ over two variables z_1 and z_2 , and the red contours represent the corresponding levels for an approximating distribution $q(\mathbf{z})$ over the same variables given by the product of two independent univariate Gaussian distributions whose parameters are obtained by minimization of (a) the Kullback-Leibler divergence $\text{KL}(q\|p)$, and (b) the reverse Kullback-Leibler divergence $\text{KL}(p\|q)$.



Mode seeking versus mode covering

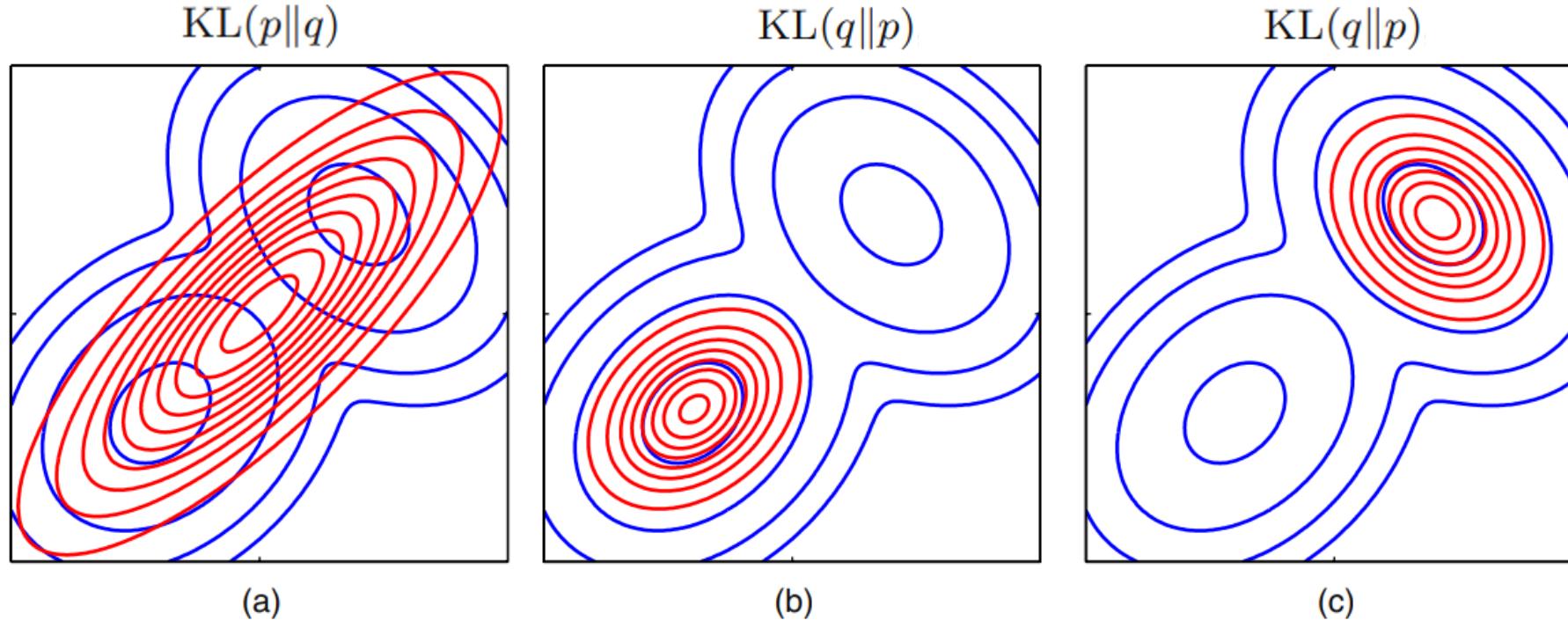


Figure 10.3 Another comparison of the two alternative forms for the Kullback-Leibler divergence. (a) The blue contours show a bimodal distribution $p(\mathbf{Z})$ given by a mixture of two Gaussians, and the red contours correspond to the single Gaussian distribution $q(\mathbf{Z})$ that best approximates $p(\mathbf{Z})$ in the sense of minimizing the Kullback-Leibler divergence $\text{KL}(p\|q)$. (b) As in (a) but now the red contours correspond to a Gaussian distribution $q(\mathbf{Z})$ found by numerical minimization of the Kullback-Leibler divergence $\text{KL}(q\|p)$. (c) As in (b) but showing a different local minimum of the Kullback-Leibler divergence.

α -divergences

Black-box α -Divergence minimization, Hernandez-Lobato et al, 2015

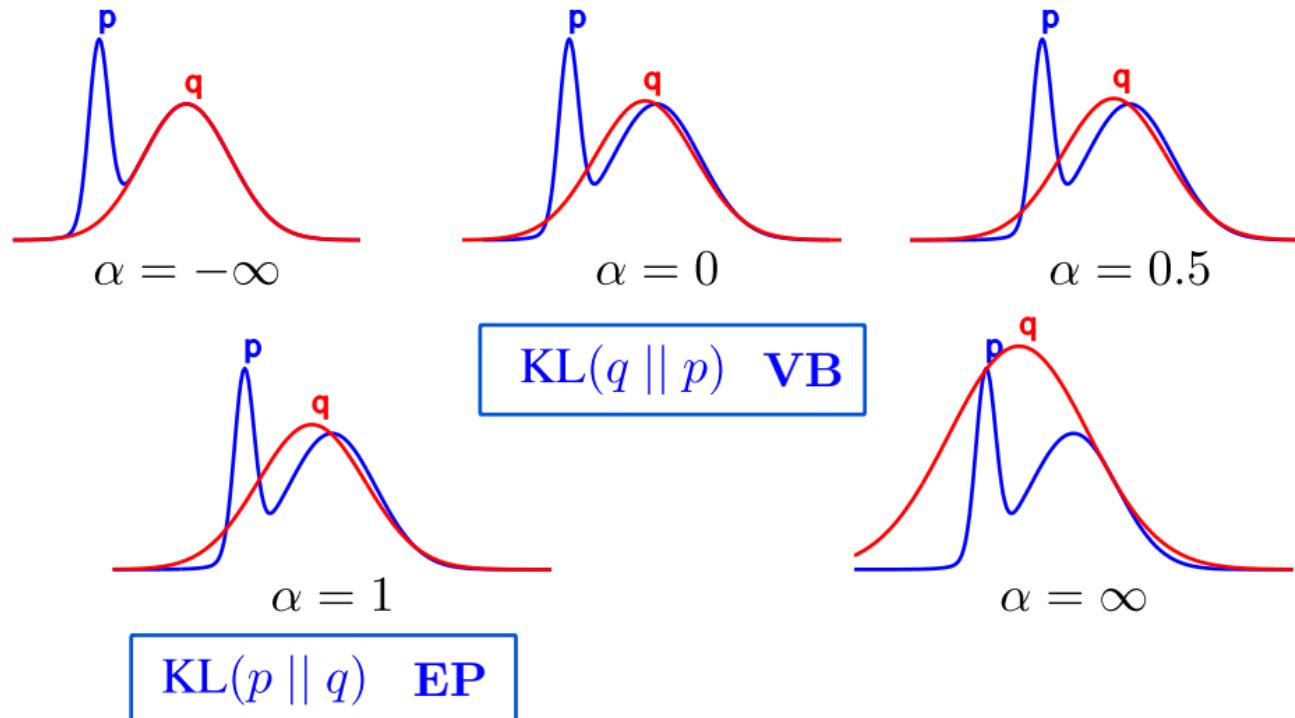


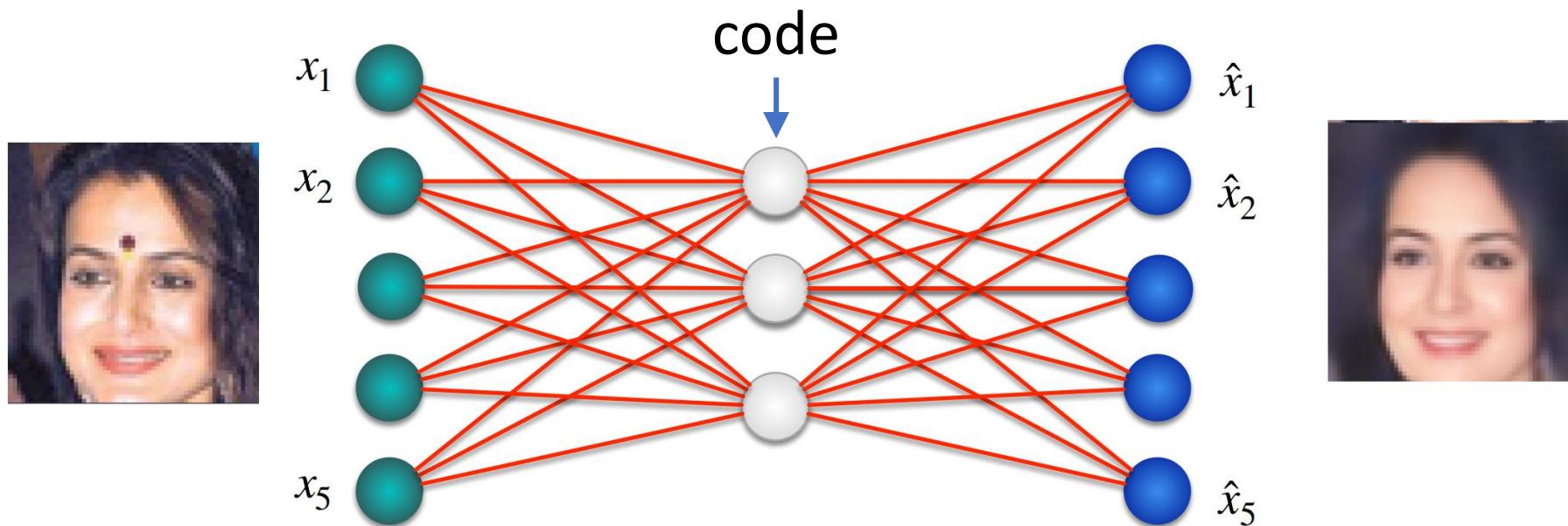
Figure 2. An illustration of approximating distributions by α -divergence minimization. Here p and q shown in the graphs are unnormalized probability densities. Reproduced from [Minka \(2005\)](#). Best viewed in color.

Outline

1. Variational auto-encoders
 - a. Variational inference
 - b. **Variational auto-encoders**
2. Normalizing flows
 - a. Normalizing flows for variational inference
 - b. Generative normalizing flows
3. Denoising diffusion models
4. Generative models for lossless compression
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Auto-encoders

Copy input to output while going through a bottleneck

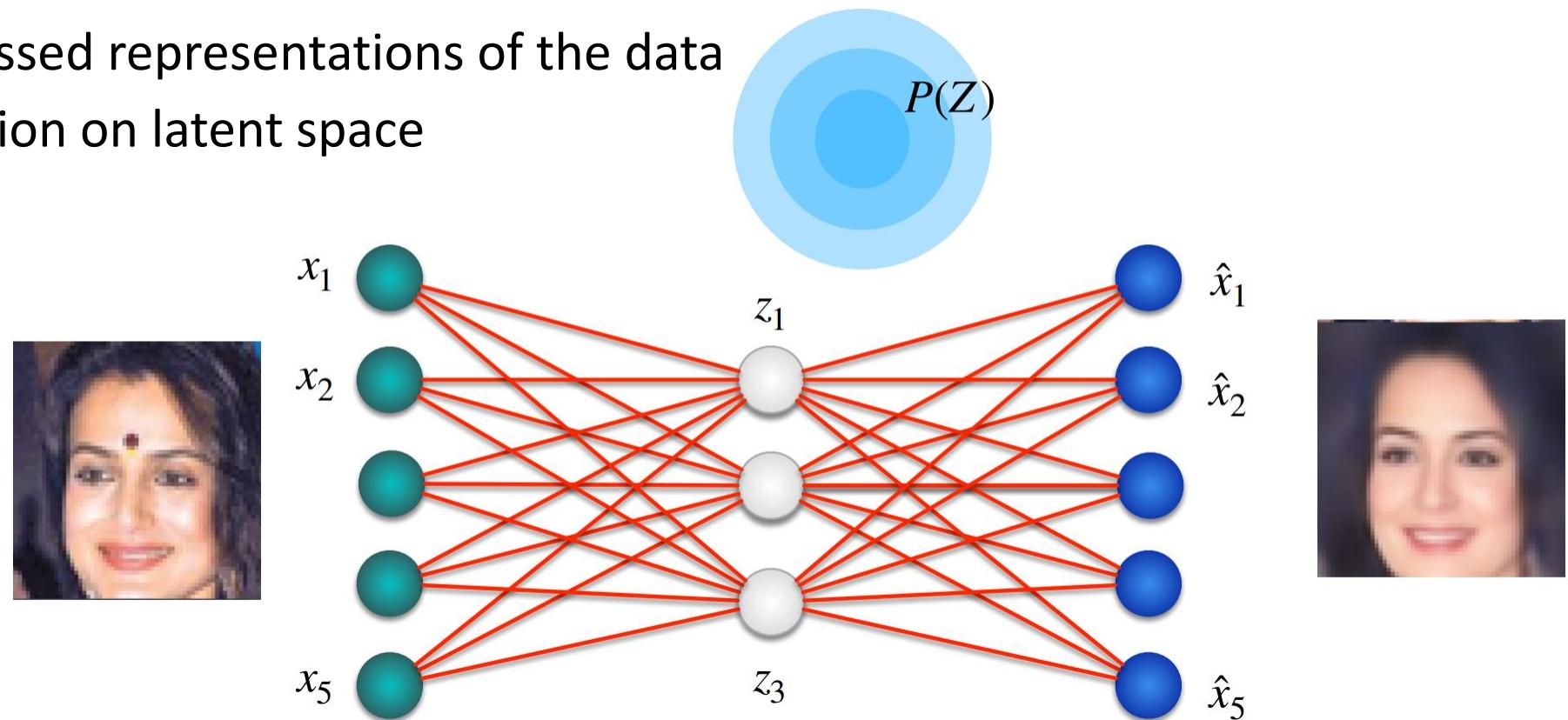


The code in the bottleneck should be a compressed representation of data x

Auto-encoders as generative models

Generative auto-encoders:

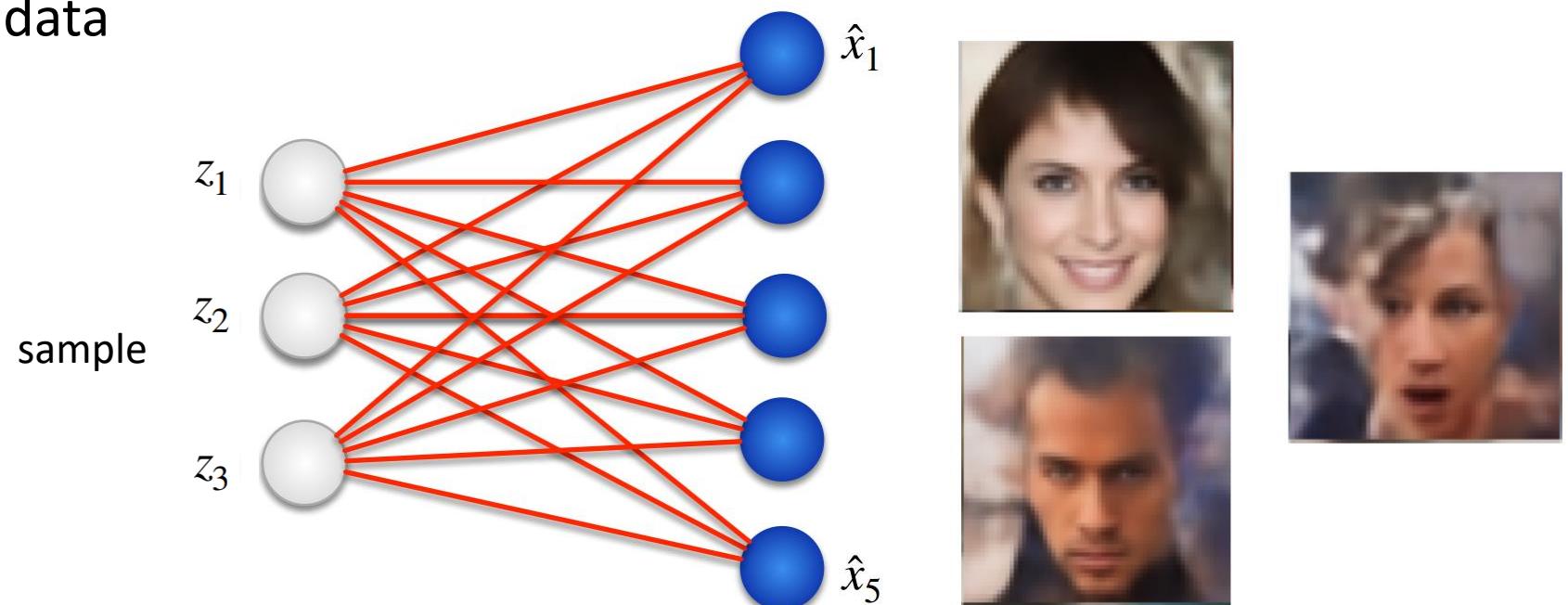
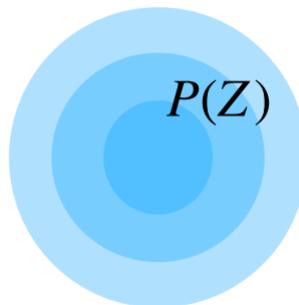
- Learn compressed representations of the data
- prior distribution on latent space



Auto-encoders as generative models

Generative auto-encoders:

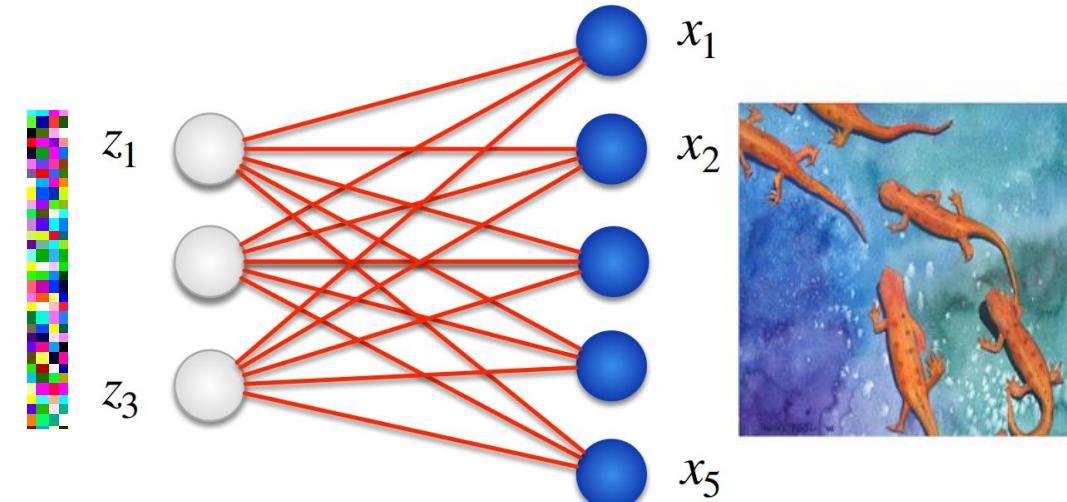
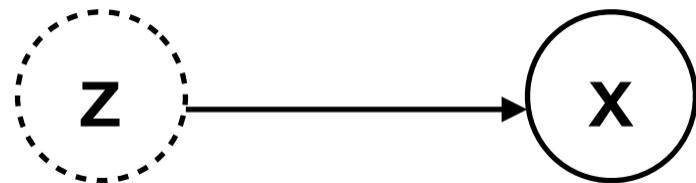
- Learn compressed representations of the data
- prior distribution on latent space
- Can generate new data



Variational auto-encoders

Kingma & Welling 2014; Rezende, Mohamed, Wierstra, 2014

Joint model $p_{\theta}(x, z) = p_{\theta}(x|z)p_{\theta}(z)$



Desired objective: $\log p_{\theta}(x) = \log \mathbb{E}_{z \sim p(z)}[p_{\theta}(x|z)]$

Surrogate objective:

$$\log p_{\theta}(x) \geq \mathbb{E}_{q(z|x)}[\log p_{\theta}(x|z)] - KL[q(z|x)||p_{\theta}(z)] = ELBO(\theta)$$

$q(z|x)$: encoding distribution

Amortized variational inference

Surrogate objective $\arg \max_{q \in Q, \theta} \mathbb{E}_{q(z|x)} [\log p_\theta(x|z)] - KL[q(z|x) || p_\theta(z)]$

Q: family of approximate posteriors.

Variational inference: optimize parameters of $q(z|x)$ for each x separately

Amortized variational inference: model parameters of $q(z|x)$ with a neural network $nn_\theta(x)$
 $q(z|x) = q_\phi(z|x)$

$$\arg \max_{\phi, \theta} \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] - KL[q_\phi(z|x) || p_\theta(z)]$$

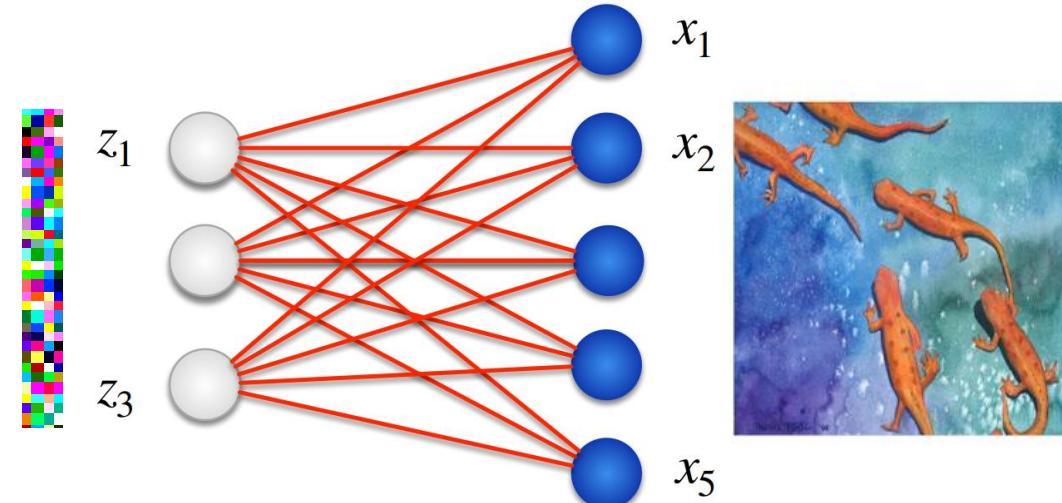
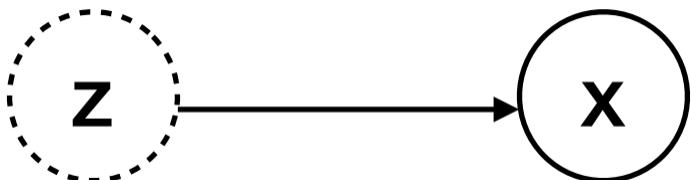
Frequently used approximate posterior: $q_\phi(z|x) = \mathcal{N}(z|\mu_\phi(x), \text{diag}(\sigma_\phi^2(x)))$

Frequently used prior: $p_\theta(z) = p(z) = \mathcal{N}(z|0, I)$

Variational auto-encoders

Kingma & Welling 2014; Rezende, Mohamed, Wierstra, 2014

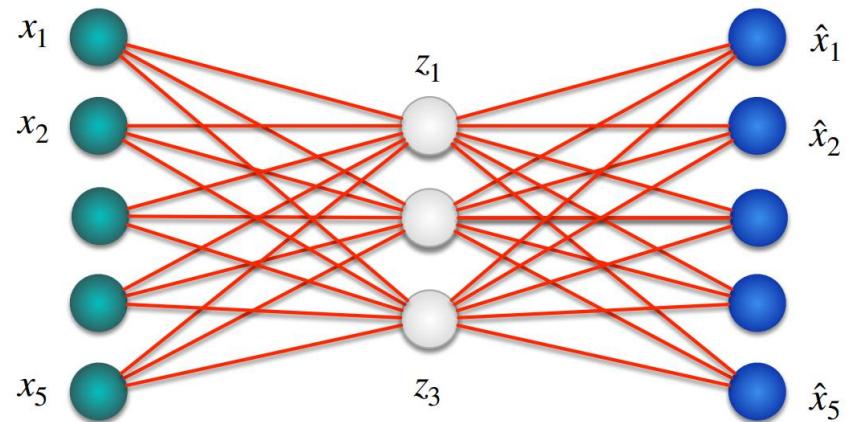
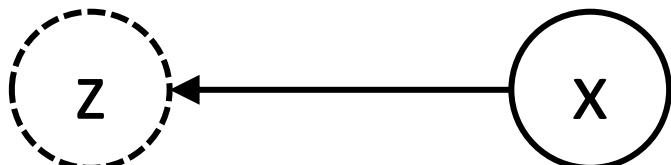
Joint model $p_{\theta}(x, z) = p_{\theta}(x|z)p_{\theta}(z)$



Surrogate objective:

$$\log p_{\theta}(x) \geq \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - KL[q_{\phi}(z|x)||p_{\theta}(z)] = ELBO(\theta, \phi)$$

$q(z|x)$: encoding distribution



VAE objective: two for one

An introduction to variational auto-encoders, Kingma & Welling, 2019

$$\arg \max_{q \in Q, \theta} \mathbb{E}_{q(z|x)} [\log p_\theta(x|z)] - KL[q_\phi(z|x) || p_\theta(z)]$$

$$\log p_\theta(x) = KL(q_\phi(z|x) || p_\theta(z|x)) + ELBO(\theta, \phi)$$

Maximizing the ELBO w.r.t. parameters θ and ϕ leads to optimization of two quantities:

1. $\log p_\theta(x)$ does not depend on ϕ so maximizing $ELBO(\theta, \phi)$ wrt ϕ makes $q_\phi(z|x)$ a better approximation to the intractable $p_\theta(z|x)$
2. We are approximately maximizing $\log p_\theta(x) \rightarrow$ improves generative model.

Optimizing the parameters of the generative model

An introduction to variational auto-encoders, Kingma & Welling, 2019

Unbiased gradients of the ELBO w.r.t. the generative model parameters θ are simple to obtain:

$$\nabla_{\theta} \mathcal{L}_{\theta, \phi}(\mathbf{x}) = \nabla_{\theta} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})] \quad (2.14)$$

$$= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x}))] \quad (2.15)$$

$$\simeq \nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z}) - \log q_{\phi}(\mathbf{z}|\mathbf{x})) \quad (2.16)$$

$$= \nabla_{\theta} (\log p_{\theta}(\mathbf{x}, \mathbf{z})) \quad (2.17)$$

The last line (eq. (2.17)) is a simple Monte Carlo estimator of the second line (eq. (2.15)), where \mathbf{z} in the last two lines (eq. (2.16) and eq. (2.17)) is a random sample from $q_{\phi}(\mathbf{z}|\mathbf{x})$.

Optimizing the parameters of the encoder

An introduction to variational auto-encoders, Kingma & Welling, 2019

Unbiased gradients w.r.t. the *variational* parameters ϕ are more difficult to obtain, since the ELBO's expectation is taken w.r.t. the distribution $q_\phi(\mathbf{z}|\mathbf{x})$, which is a function of ϕ . I.e., in general:

$$\nabla_\phi \mathcal{L}_{\theta,\phi}(\mathbf{x}) = \nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \quad (2.18)$$

$$\neq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\nabla_\phi (\log p_\theta(\mathbf{x}, \mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}))] \quad (2.19)$$

Reparameterization: separating randomness & parameters

Kingma & Welling 2014; Rezende, Mohamed, Wierstra, 2014

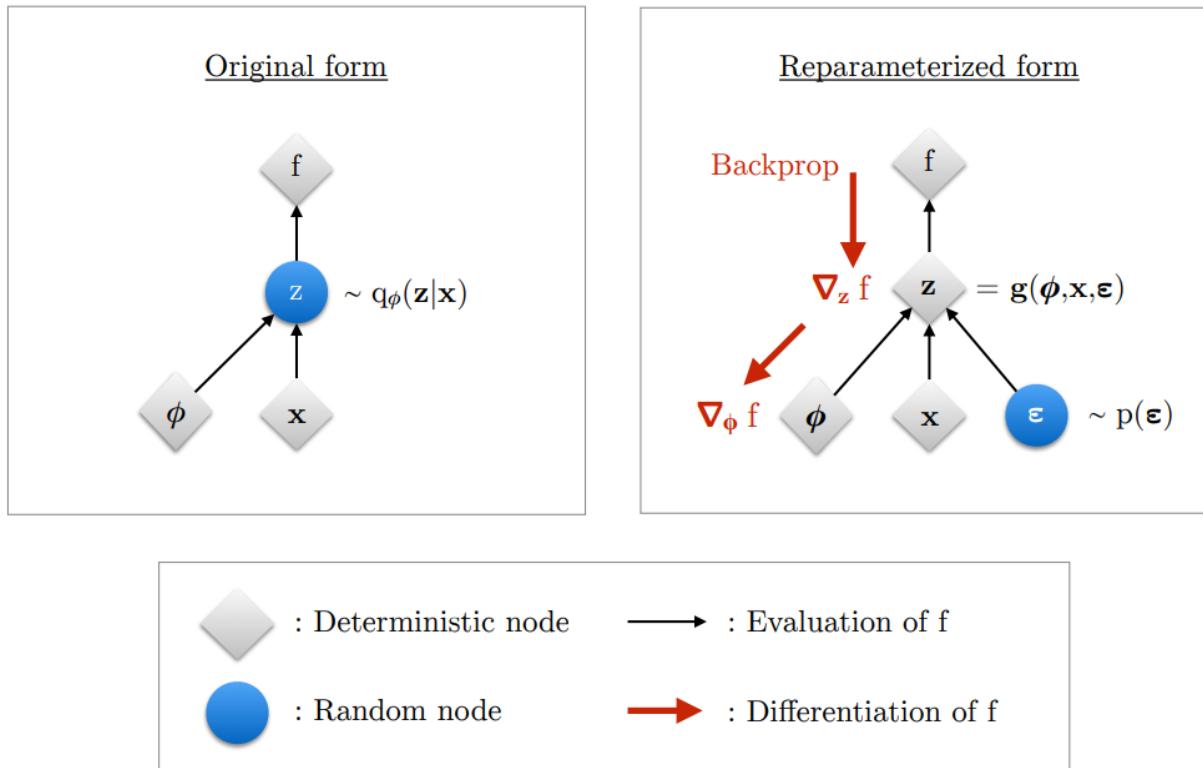
Rewrite $z \sim q_\phi(z|x)$ as a transformation of another random variable ϵ that does not depend ϕ and x :

$$z = g(\epsilon, \phi, x) \quad \epsilon \sim p(\epsilon)$$

Then for any $f(z)$ we can do the following

$$\begin{aligned} \nabla_\phi \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [f(\mathbf{z})] &= \nabla_\phi \mathbb{E}_{p(\epsilon)} [f(\mathbf{z})] \\ &= \mathbb{E}_{p(\epsilon)} [\nabla_\phi f(\mathbf{z})] \\ &\simeq \nabla_\phi f(\mathbf{z}) \end{aligned}$$

Reparameterization trick



Reparameterizing the ELBO

$$\begin{aligned} \text{ELBO}(\theta, \phi) &= \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x, z) - \log q_\phi(z|x)] \\ &= \mathbb{E}_{p(\epsilon)} [\log p_\theta(x, z(\epsilon, \phi, x)) - \log q_\phi(z(\epsilon, \phi, x)|x)] \end{aligned}$$

$$z = g(\epsilon, \phi, x) \quad \epsilon \sim p(\epsilon)$$

Gradients wrt encoder parameters ϕ can be computed with automatic backprop

$$\nabla_\phi \text{ELBO}(\theta, \phi) = \mathbb{E}_{p(\epsilon)} [\nabla_\phi \log p_\theta(x, z(\epsilon, \phi, x)) - \nabla_\phi \log q_\phi(z(\epsilon, \phi, x)|x)]$$

Reparameterization for factorized Gaussian distributions

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \prod_i q_{\phi}(z_i|\mathbf{x}) = \prod_i \mathcal{N}(z_i; \mu_i, \sigma_i^2)$$

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

Reparameterization: $\epsilon \sim \mathcal{N}(0, \mathbf{I})$

$$(\boldsymbol{\mu}, \log \boldsymbol{\sigma}) = \text{EncoderNeuralNet}_{\phi}(\mathbf{x})$$

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \epsilon$$

Reparameterization for full-covariance Gaussian distributions

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

Reparameterization:

$$\begin{aligned}\epsilon &\sim \mathcal{N}(0, \mathbf{I}) \\ \mathbf{z} &= \boldsymbol{\mu} + \mathbf{L}\epsilon\end{aligned}$$

Cholesky decomposition of the covariance matrix: $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$
(\mathbf{L} is lower triangular and has nonzero diagonal entries)

Results from Kingma & Welling 2014

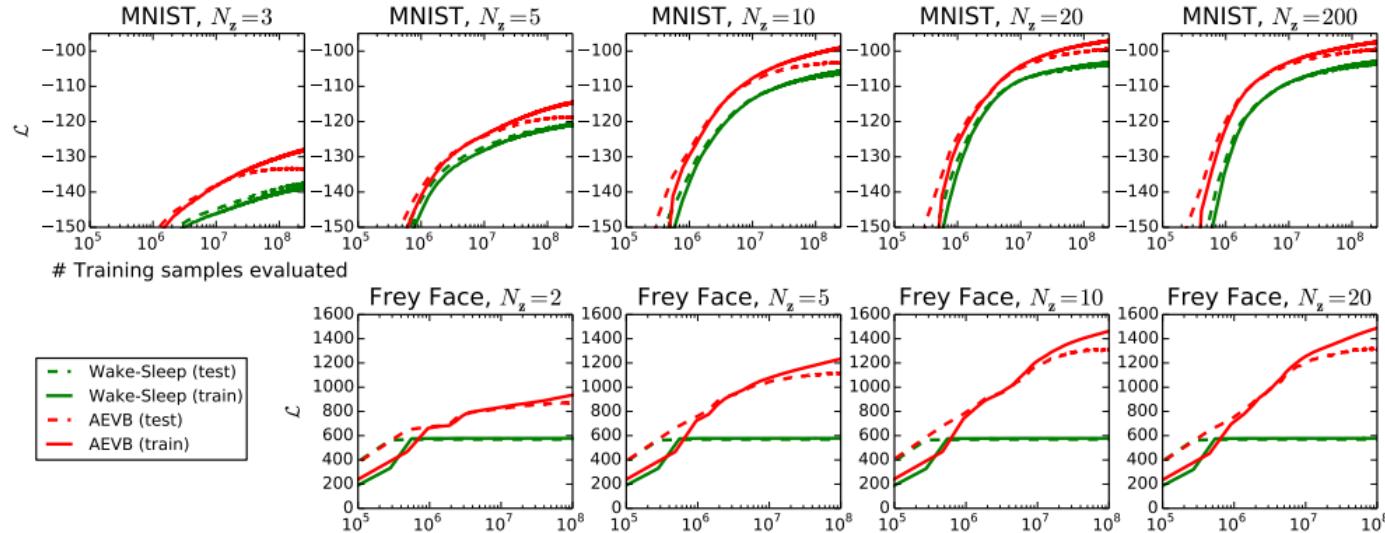
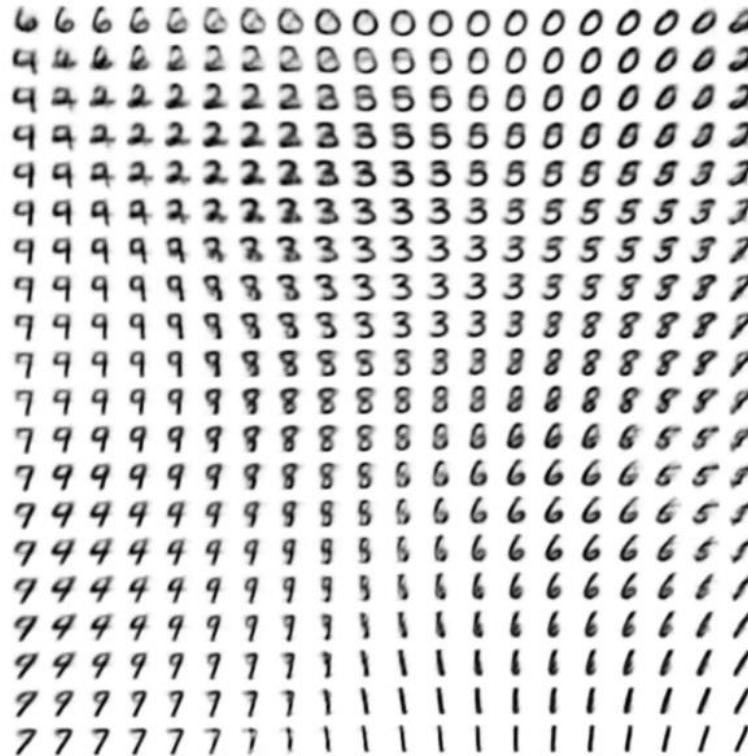


Figure 2: Comparison of our AEVB method to the wake-sleep algorithm, in terms of optimizing the lower bound, for different dimensionality of latent space (N_z). Our method converged considerably faster and reached a better solution in all experiments. Interestingly enough, more latent variables does not result in more overfitting, which is explained by the regularizing effect of the lower bound. Vertical axis: the estimated average variational lower bound per datapoint. The estimator variance was small (< 1) and omitted. Horizontal axis: amount of training points evaluated. Computation took around 20-40 minutes per million training samples with a Intel Xeon CPU running at an effective 40 GFLOPS.

Results from Kingma & Welling 2014



(a) Learned Frey Face manifold



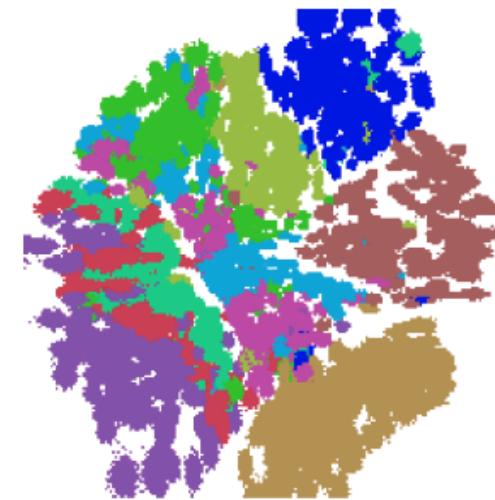
(b) Learned MNIST manifold

Figure 4: Visualisations of learned data manifold for generative models with two-dimensional latent space, learned with AEVB. Since the prior of the latent space is Gaussian, linearly spaced coordinates on the unit square were transformed through the inverse CDF of the Gaussian to produce values of the latent variables \mathbf{z} . For each of these values \mathbf{z} , we plotted the corresponding generative $p_{\theta}(\mathbf{x}|\mathbf{z})$ with the learned parameters θ .

Results from Rezende et al. 2014



(a) Left: Training data. Middle: Sampled pixel probabilities. Right: Model samples



(b) 2D embedding.

Figure 3. Performance on the MNIST dataset. For the visualisation, each colour corresponds to one of the digit classes.

Marginal likelihood estimation: importance sampling

Rezende et al. 2014, Burda et al. 2016

$$ELBO(\theta, \phi) = \mathbb{E}_{q_\phi(z|x)} [\log \frac{p_\theta(x, z)}{q_\phi(z|x)}] \approx \frac{1}{K} \sum_{k=1}^K \log \frac{p_\theta(x, z^{(k)})}{q_\phi(z^{(k)}|x)}$$

Tighter bound is obtained by importance sampling:

$$ELBO_{IW}^K(\theta, \phi) = \mathbb{E}_{z^{(1)}, \dots, z^{(K)} \sim q_\phi(z|x)} [\log (\frac{1}{K} \sum_{k=1}^K \frac{p_\theta(x, z^{(k)})}{q_\phi(z^{(k)}|x)})] \leq \log p_\theta(x)$$

A very incomplete selection of advances in VAEs: Ladder VAEs

Ladder Variational autoencoders, Sonderby et al. 2016

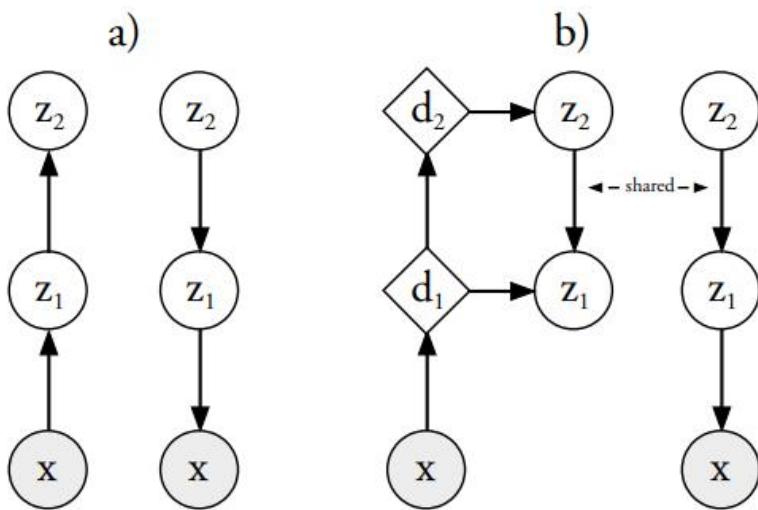


Figure 1: Inference (or encoder/recognition) and generative (or decoder) models for a) VAE and b) LVAE. Circles are stochastic variables and diamonds are deterministic variables.

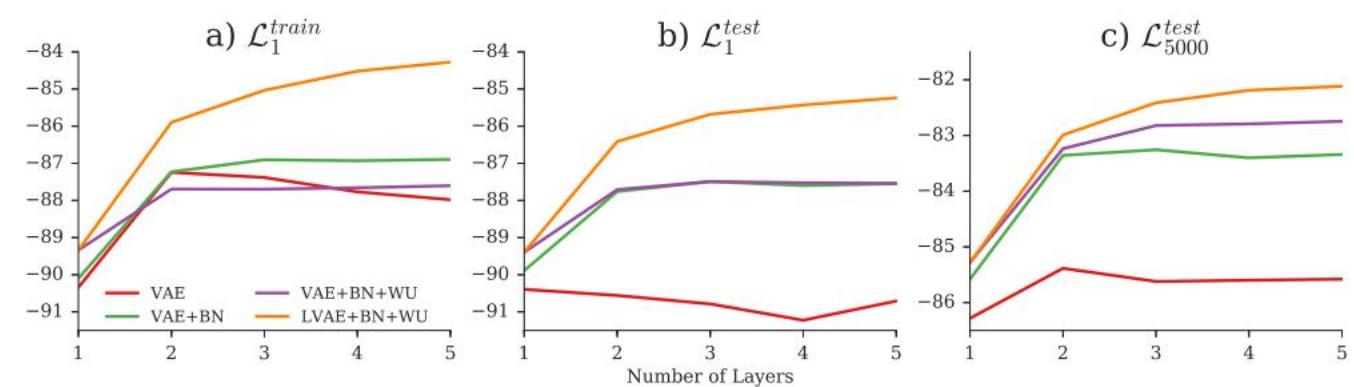


Figure 3: MNIST log-likelihood values for VAEs and the LVAE model with different number of latent layers, Batch normalization (BN) and Warm-up (WU). a) Train log-likelihood, b) test log-likelihood and c) test log-likelihood with 5000 importance samples.

Follow up work: Biva: a very deep hierarchy of latent variables for generative modeling, Maaloe et al. 2019

A very incomplete selection of advances in VAEs: NVAE

NVAE: a deep hierarchical variational autoencoder, 2020

In summary, we make the following contributions: i) We propose a novel deep hierarchical VAE, called NVAE, with depthwise convolutions in its generative model. ii) We propose a new residual parameterization of the approximate posteriors. iii) We stabilize training deep VAEs with spectral regularization. iv) We provide practical solutions to reduce the memory burden of VAEs. v) We show that deep hierarchical VAEs can obtain state-of-the-art results on several image datasets, and can produce high-quality samples even when trained with the original VAE objective. To the best of our knowledge, NVAE is the first successful application of VAEs to images as large as 256×256 pixels.

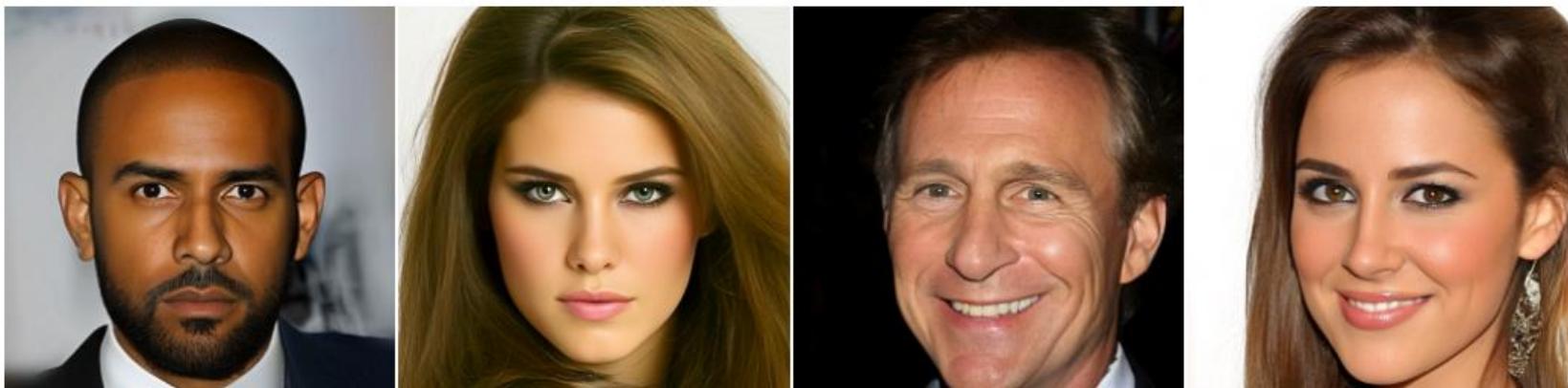


Figure 1: 256×256 -pixel samples generated by NVAE, trained on CelebA HQ [28].

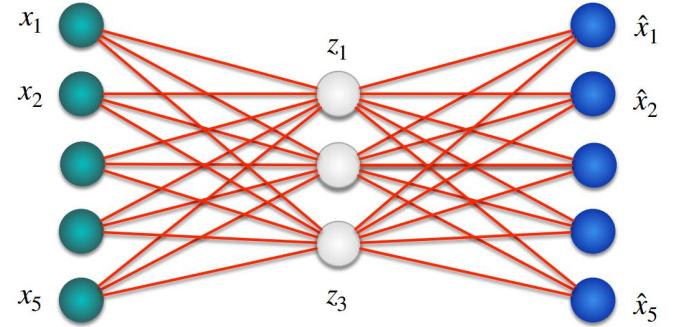
Outline

1. Variational auto-encoders
 - a. Variational inference
 - b. Variational auto-encoders
2. **Normalizing flows**
 - a. Normalizing flows for variational inference
 - b. Generative normalizing flows
3. Denoising diffusion models
4. Generative models for lossless compression
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Variational auto-encoders

Surrogate objective:

$$\log p_\theta(x) \geq \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - KL[q_\phi(z|x)||p_\theta(z)] = ELBO(\theta, \phi)$$

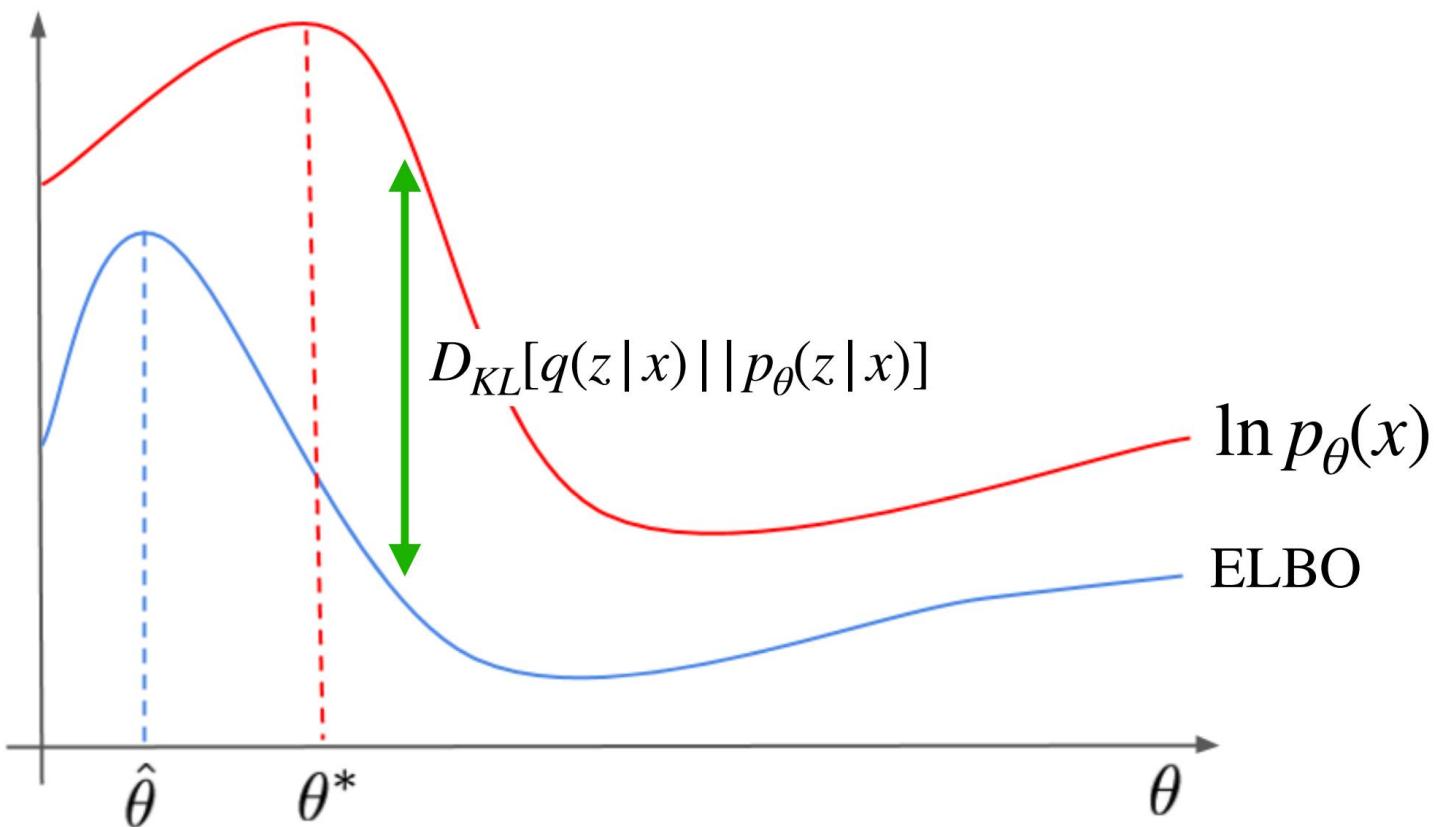


Equivalent to:

$$ELBO(\theta, \phi) = \log p_\theta(x) - KL(q_\phi(z|x)||p_\theta(z|x))$$

- $q_\phi(z|x)$ tries to approximate $p_\theta(z|x)$
- Bad approximations lead to a large gap!

Approximate posteriors



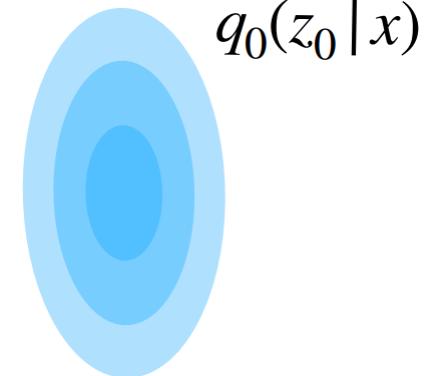
Looser lower bounds do not share the same local maxima as the log likelihood

Variational inference with normalizing flows

Rezende & Mohamed, ICML 2015

Sample from simple distribution $z_0 \in \mathbb{R}^D$

$$z_0 \sim q_0(z|x) = \mathcal{N}(z|\mu(x), \text{diag}(\sigma^2(x)))$$

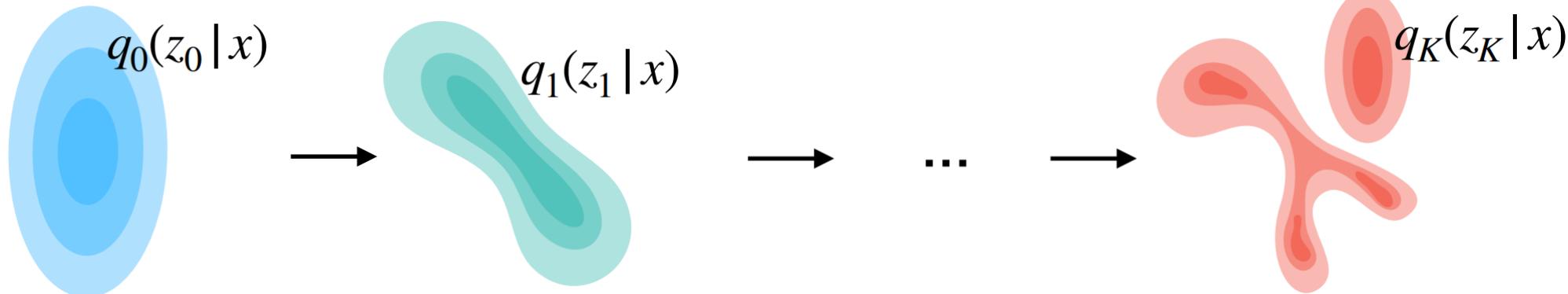


Apply a sequence of invertible transformations

$$z_K = f_K \circ \dots \circ f_2 \circ f_1(z_0)$$

$$z_0 \rightarrow z_1 \rightarrow \dots \rightarrow z_K$$

$$f_k : \mathbb{R}^D \mapsto \mathbb{R}^D$$

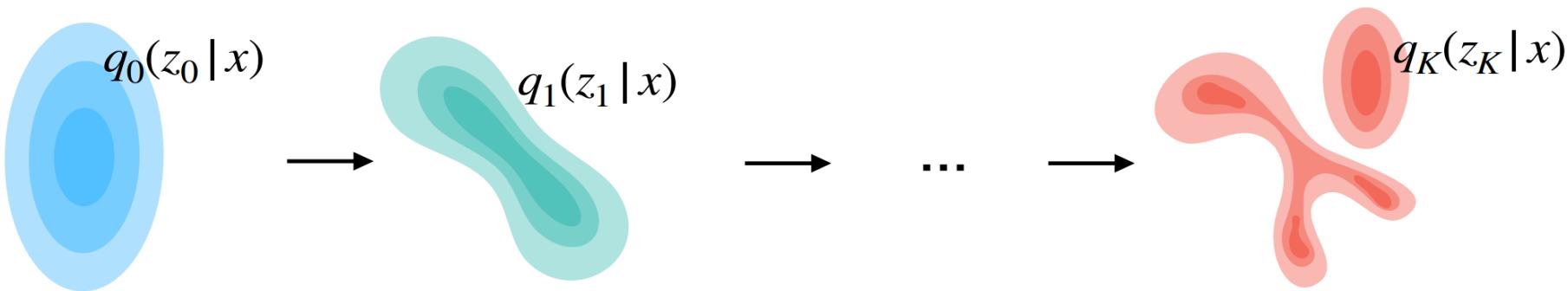


Variational inference with normalizing flows

Rezende & Mohamed, ICML 2015

Apply a sequence of invertible transformations $f_k : \mathbb{R}^D \mapsto \mathbb{R}^D$

$$z_K = f_K \circ \dots \circ f_2 \circ f_1(z_0) \quad z_0 \rightarrow z_1 \rightarrow \dots \rightarrow z_K$$



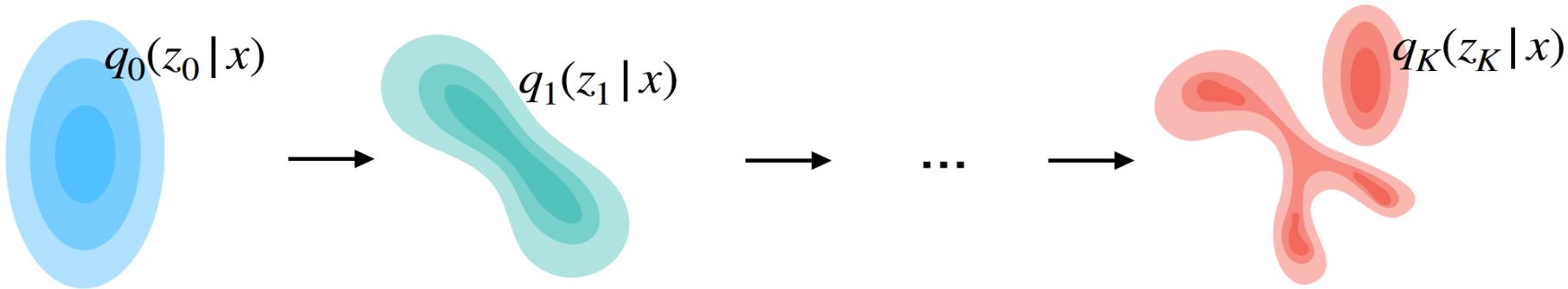
For each transformation $z_k = f_k(z_{k-1})$

$$q_k(z_k) = q_{k-1}(z_{k-1}) \left| \det \frac{\partial f_k}{\partial z_{k-1}} \right|^{-1} \quad \rightarrow \quad \ln q_K(z_K) = \ln q_0(z_0) - \sum_{k=1}^K \ln \left| \det \frac{\partial f_k}{\partial z_{k-1}} \right|$$

Practical normalizing flows for variational inference

1. Flexible invertible transformations

Inflexible transformations lead to long sequences of flows for flexible posteriors



2. Easily computable Jacobian determinants:

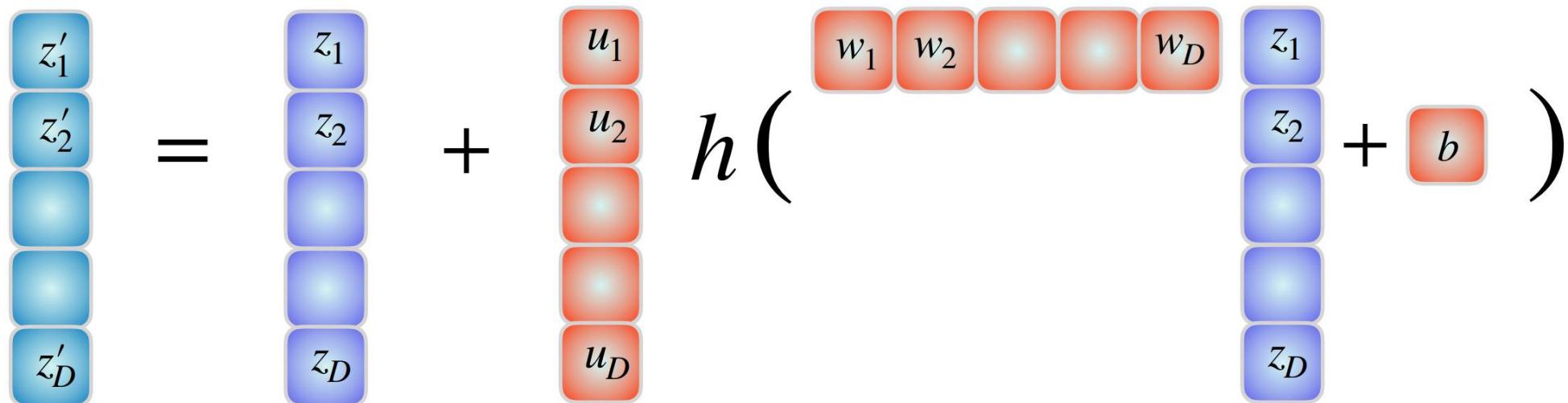
For a $D \times D$ matrix in general $O(D^3)$

Planar flows

Rezende & Mohamed, ICML 2015

$$z' = z + u \ h(w^T z + b)$$

: learnable parameter

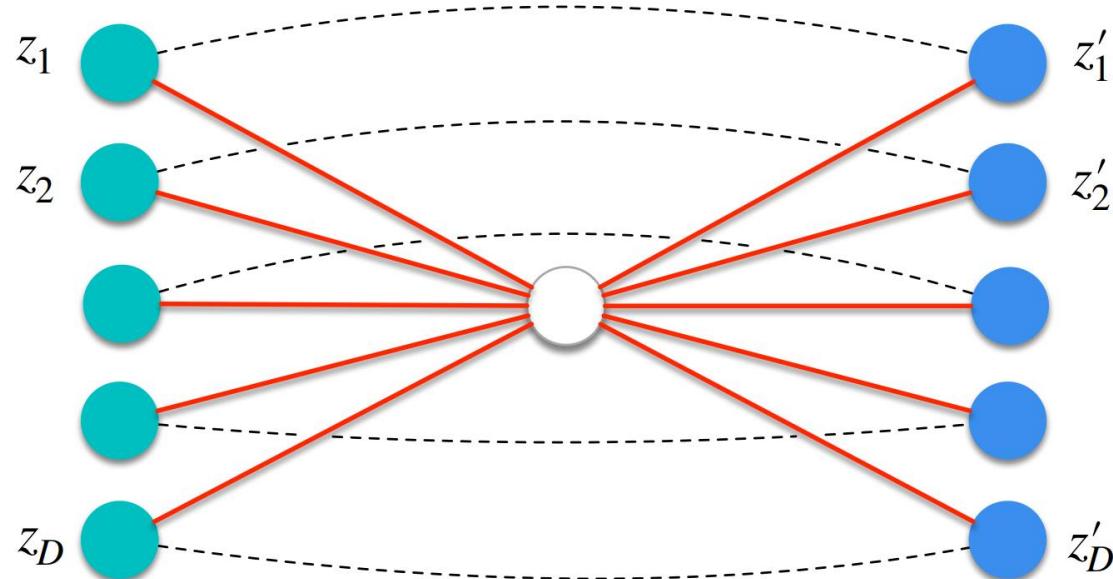


Planar flows

Rezende & Mohamed, ICML 2015

$$z' = z + u \ h(w^T z + b)$$

2 layer block single hidden
unit with a residual
connection



Planar flows: very simple jacobian determinant

$$z' = z + u h(w^T z + b)$$

$$h(z) = \tanh(z)$$

$$w^T u \geq -1$$

Simple determinant:

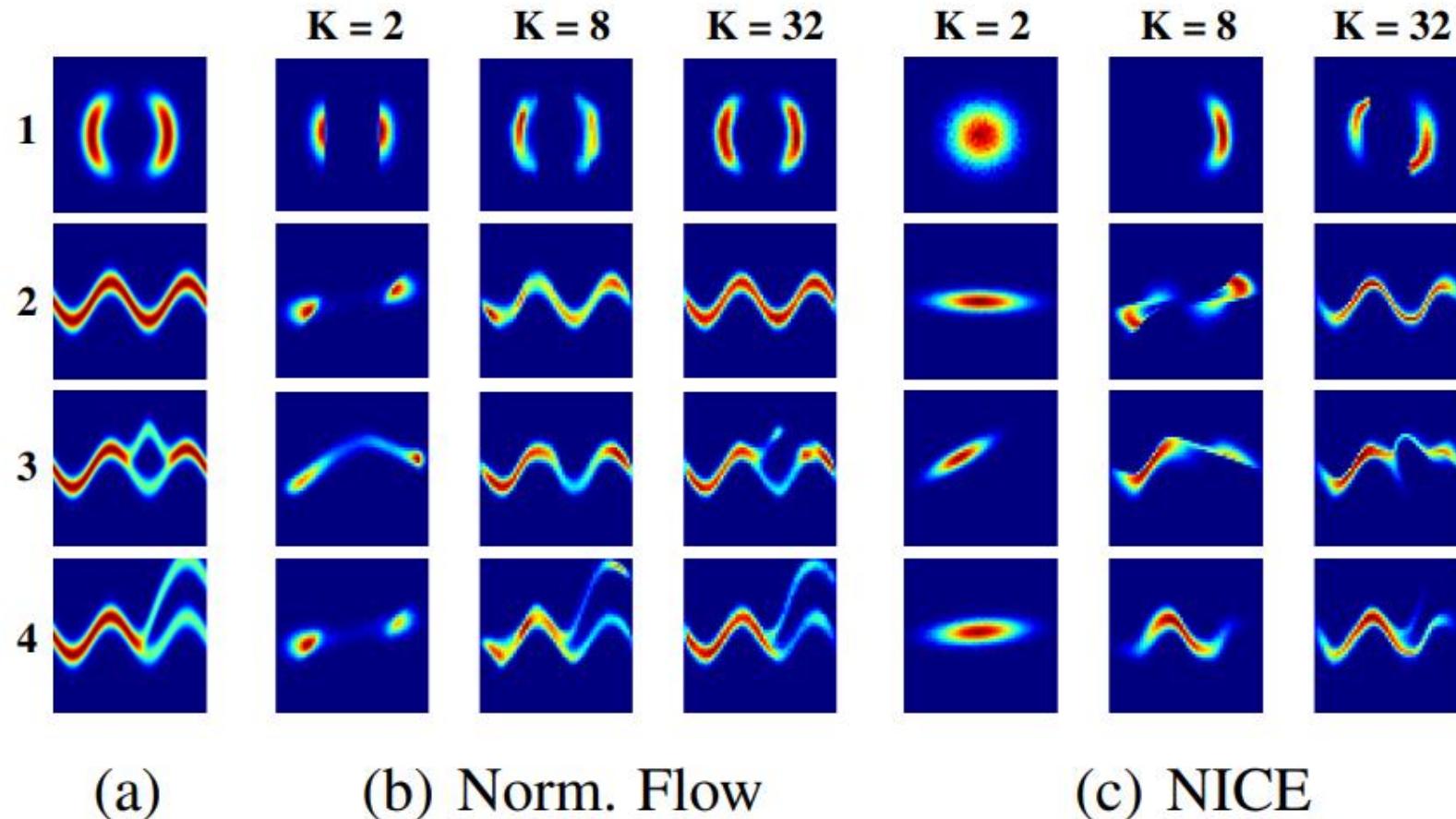
$$\begin{aligned} \det \frac{\partial z'}{\partial z} &= \det (\mathbf{1}_D + u h'(w^T z + b) w^T) \\ &= 1 + h'(w^T z + b) w^T u \end{aligned}$$

Matrix determinant lemma:

$$\det(A + uv^T) = (1 + v^T A^{-1} u) \det(A)$$

Planar flows

Rezende & Mohamed, ICML 2015

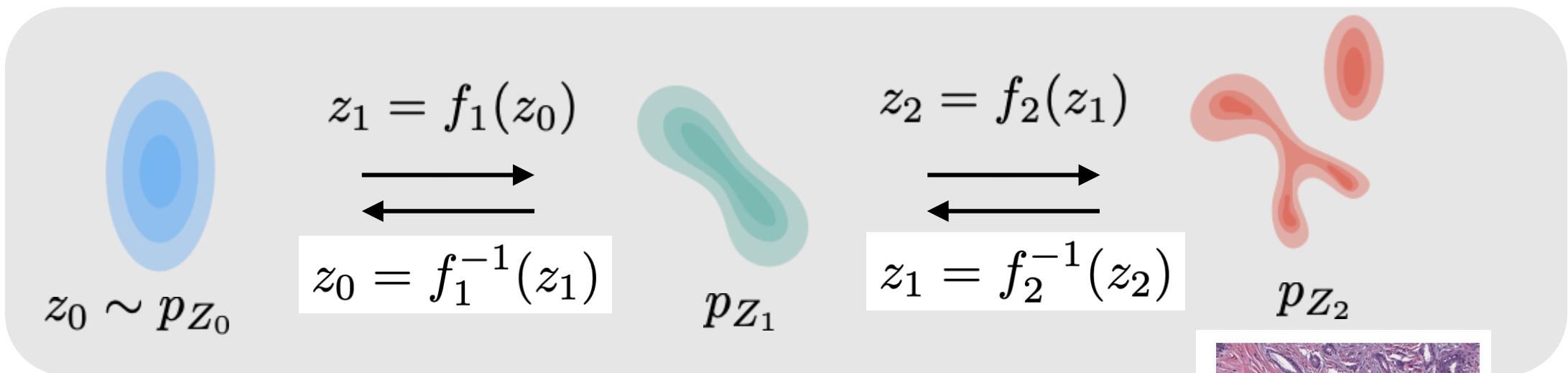


Outline

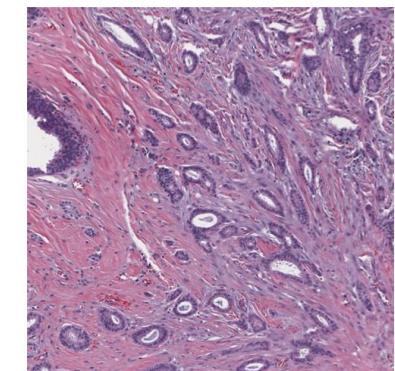
1. Variational auto-encoders
 - a. Variational inference
 - b. Variational auto-encoders
2. Normalizing flows
 - a. Normalizing flows for variational inference
 - b. **Generative normalizing flows**
3. Denoising diffusion models
4. Generative models for lossless compression
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Generative normalizing flows

Idea: apply a sequence of invertible transformations to a random variable



Additional practical requirement: The inverse function also needs to be easy to compute!



Density estimation using Real NVP

Dinh, Sohl-Dickstein, Bengio 2016



$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} \leftarrow \begin{bmatrix} x_1 \\ s^\theta(x_1) \odot x_2 + t^\theta(x_1) \end{bmatrix}$$



Non-volume preserving



$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leftarrow \begin{bmatrix} z_1 \\ (z_2 - t^\theta(z_1)) \odot s^\theta(z_1) \end{bmatrix}$$



Image generation with RealNVP

Dinh, Sohl-Dickstein, Bengio 2016

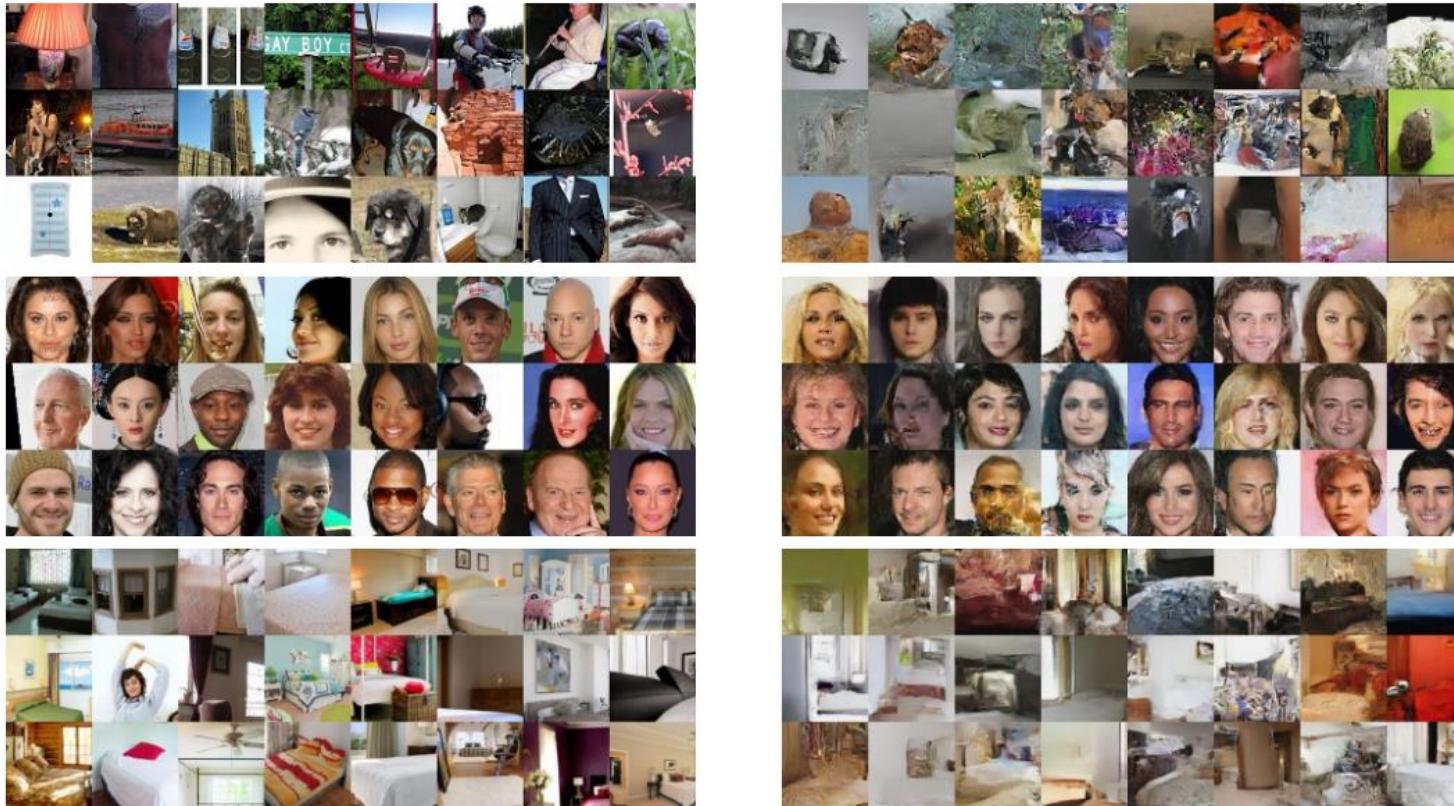


Figure 5: On the left column, examples from the dataset. On the right column, samples from the model trained on the dataset. The datasets shown in this figure are in order: CIFAR-10, Imagenet (32 × 32), Imagenet (64 × 64), CelebA, LSUN (bedroom).

Masked autoregressive flow for density estimation

Papamakarios et al. 2017



Latent z to data x: autoregressive

$$x_i = \mu_i(x_{1:i-1}) + z_i \odot \sigma_i(x_{1:i-1})$$

Data x to latent z: parallel

$$z_i = (x_i - \mu_i(x_{1:i-1}))\sigma_i^{-1}(x_{1:i-1})$$



Glow: generative flow with invertible 1x1 convolutions

Kingma & Dhariwal 2018

1x1 convolution: invertible matrix multiplication along the channel dimension



Figure 7: Samples from model trained on 5-bit LSUN bedrooms, at temperature 0.875. Resolutions 64, 96 and 128 respectively⁴

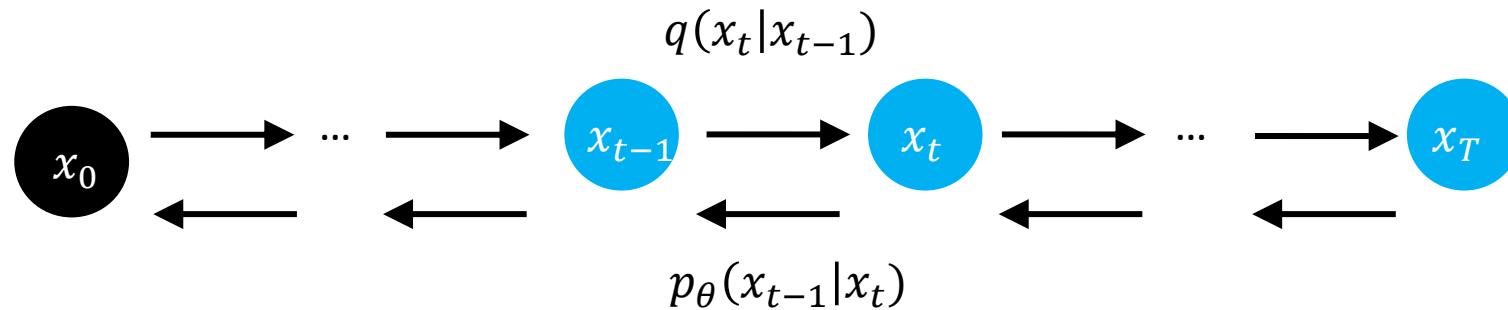
Outline

1. Variational auto-encoders
 - a. Variational inference
 - b. Variational auto-encoders
2. Normalizing flows
 - a. Normalizing flows for variational inference
 - b. Generative normalizing flows
- 3. Denoising diffusion models**
4. Generative models for lossless compression
 1. Basics of lossless compression
 2. Connecting likelihood-based generative models and lossless compression
 3. Discrete normalizing flows for lossless compression
 4. Bits-back coding: variational auto-encoders and lossless compression

Denoising diffusion probabilistic models

Sohl-Dickstein et al., ICML 2015, Ho et al., NeurIPS 2020, Song et al., ICLR 2021

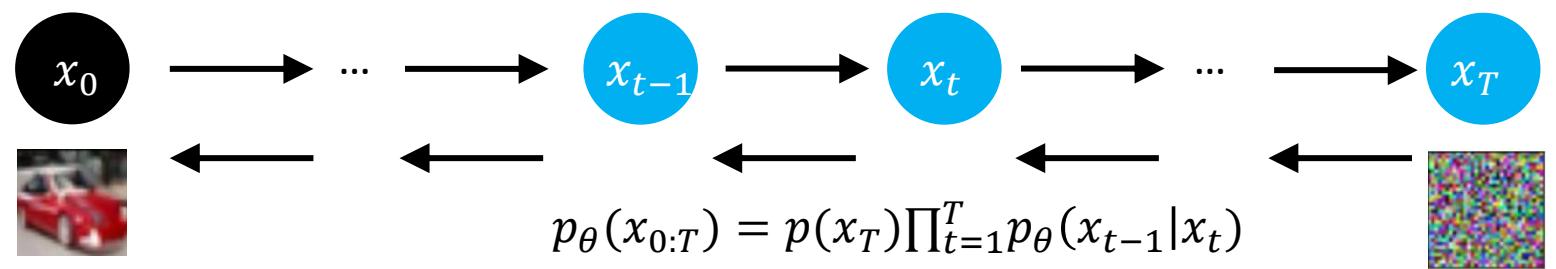
Forward process: corrupt data with noise



Reverse process: learning to denoise data



Training diffusion models



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

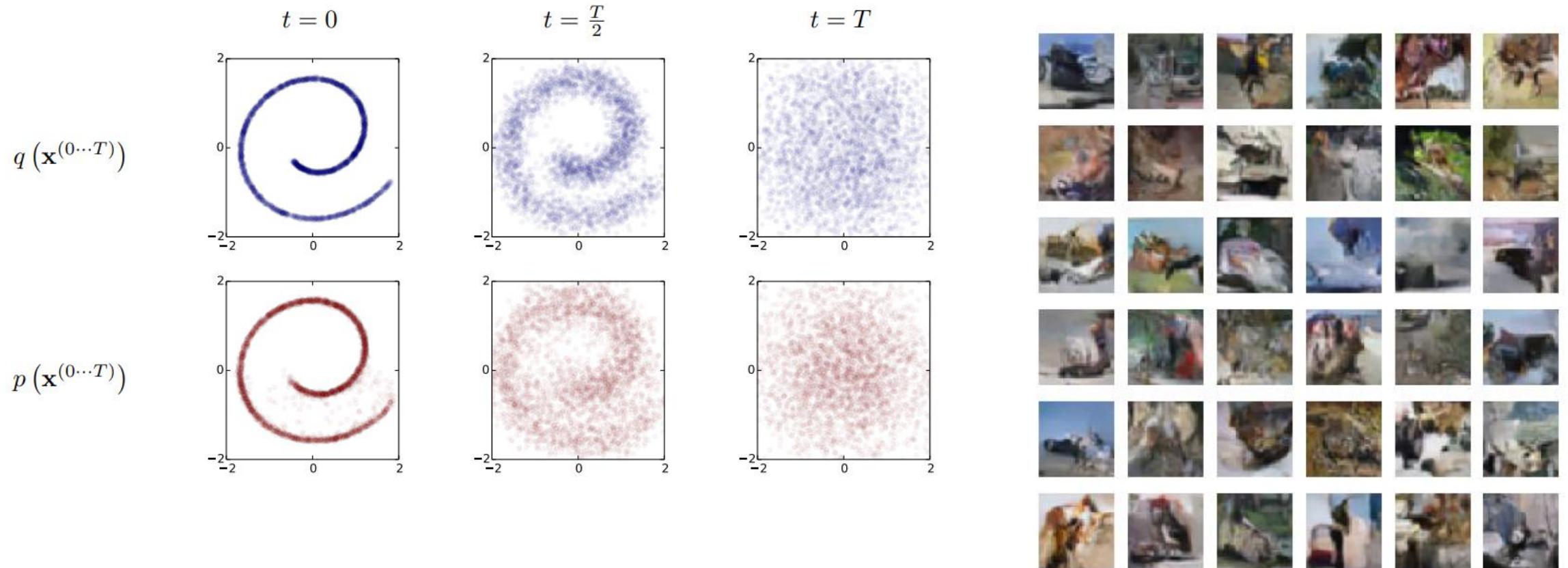
Practical requirements for $q(x_t|x_{t-1})$ to allow for efficient training of p_θ :

1. Efficient sampling of x_t from $q(x_t|x_0)$ for arbitrary time t
2. Tractable expression for $q(x_{t-1}|x_t, x_0)$.

For Gaussian or binomial $q(x_t|x_{t-1})$ (and $p_\theta(x_{t-1}|x_t)$):



Results from Sohl-Dickstein et al. 2015



Results from Ho et al. 2020

Sohl-Dickstein et al., ICML 2015, Ho et al., NeurIPS 2020, Song et al., ICLR 2021

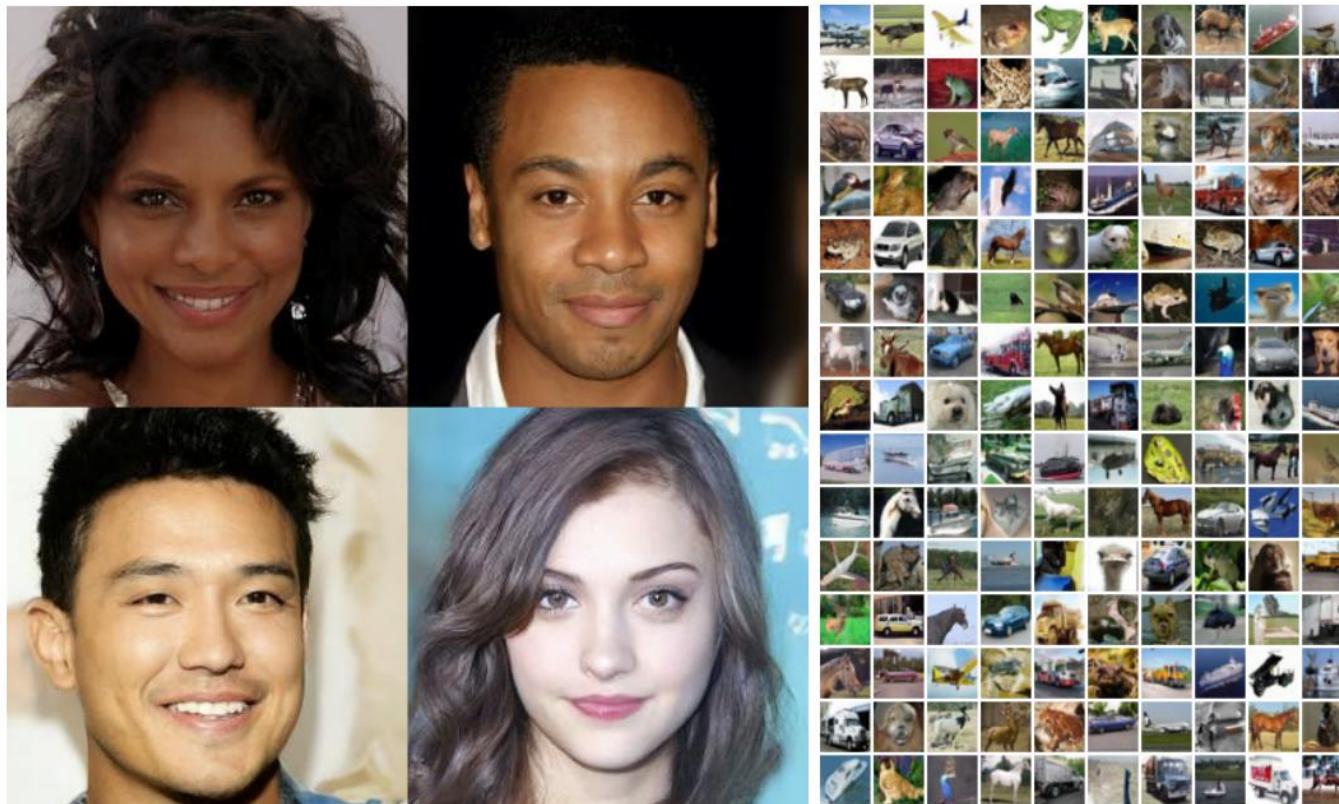


Figure 1: Generated samples on CelebA-HQ 256×256 (left) and unconditional CIFAR10 (right)

Table 1: CIFAR10 results. NLL measured in bits/dim.

Model	IS	FID	NLL Test (Train)
Conditional			
EBM [11]	8.30	37.9	
JEM [17]	8.76	38.4	
BigGAN [3]	9.22	14.73	
StyleGAN2 + ADA (v1) [29]	10.06	2.67	
Unconditional			
Diffusion (original) [53]			≤ 5.40
Gated PixelCNN [59]	4.60	65.93	3.03 (2.90)
Sparse Transformer [7]			2.80
PixelIQN [43]	5.29	49.46	
EBM [11]	6.78	38.2	
NCSNv2 [56]			31.75
NCSN [55]	8.87 ± 0.12	25.32	
SNGAN [39]	8.22 ± 0.05	21.7	
SNGAN-DDLS [4]	9.09 ± 0.10	15.42	
StyleGAN2 + ADA (v1) [29]	9.74 ± 0.05	3.26	
Ours (L , fixed isotropic Σ)	7.67 ± 0.13	13.51	≤ 3.70 (3.69)
Ours (L_{simple})	9.46 ± 0.11	3.17	≤ 3.75 (3.72)