

Introduction to Variational Auto-Encoders (VAE)

Tomasz Kajdanowicz, Piotr Bielek, Maciej Falkiewicz,
Kacper Kania, Piotr Zieliński

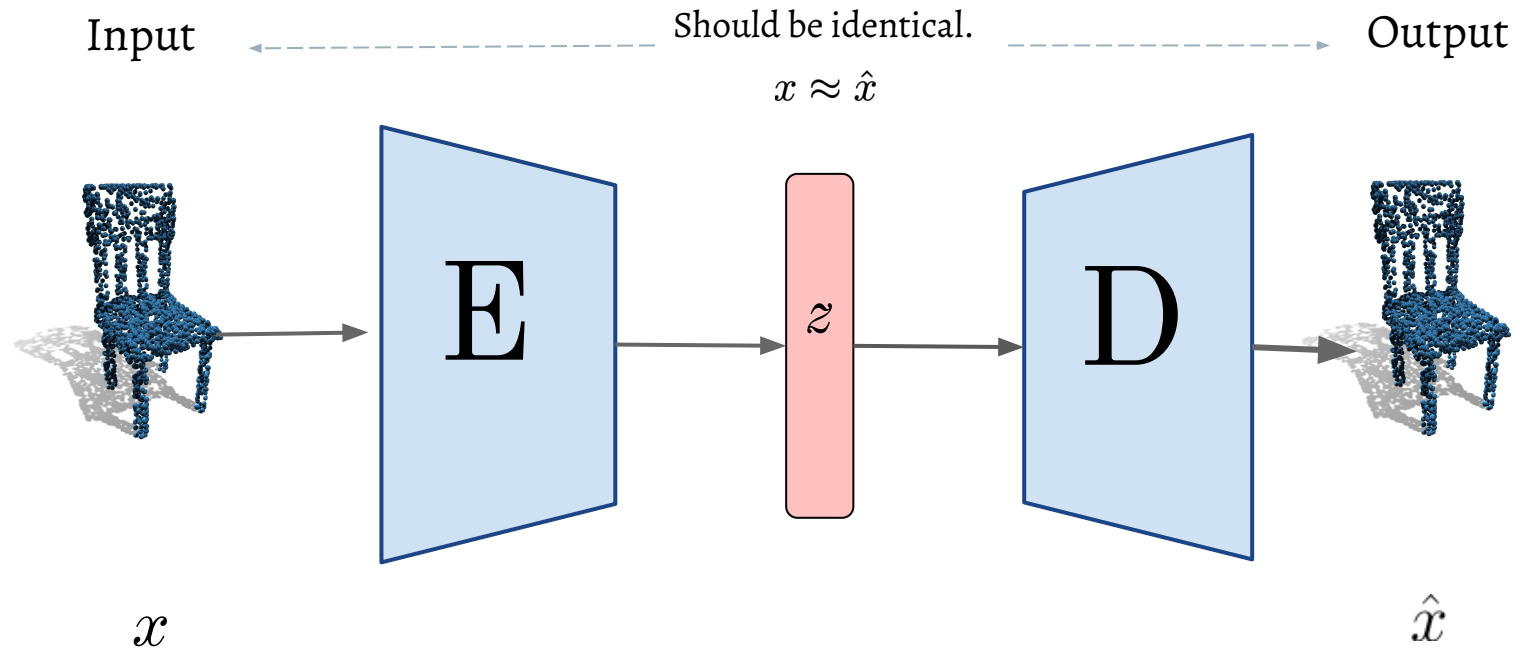
Rereferences

1. Goodfellow, Ian, "NIPS 2016 Tutorial: Generative Adversarial Networks " arXiv preprint arXiv:1701.00160 (2016)
2. Doersch, C. (2016). Tutorial on variational autoencoders. arXiv preprint arXiv:1606.05908.
3. Welling, Max; Kingma, Diederik P. (2019). "An Introduction to Variational Autoencoders". Foundations and Trends in Machine Learning. 12 (4): 307–392. arXiv:1906.02691

Further reading:

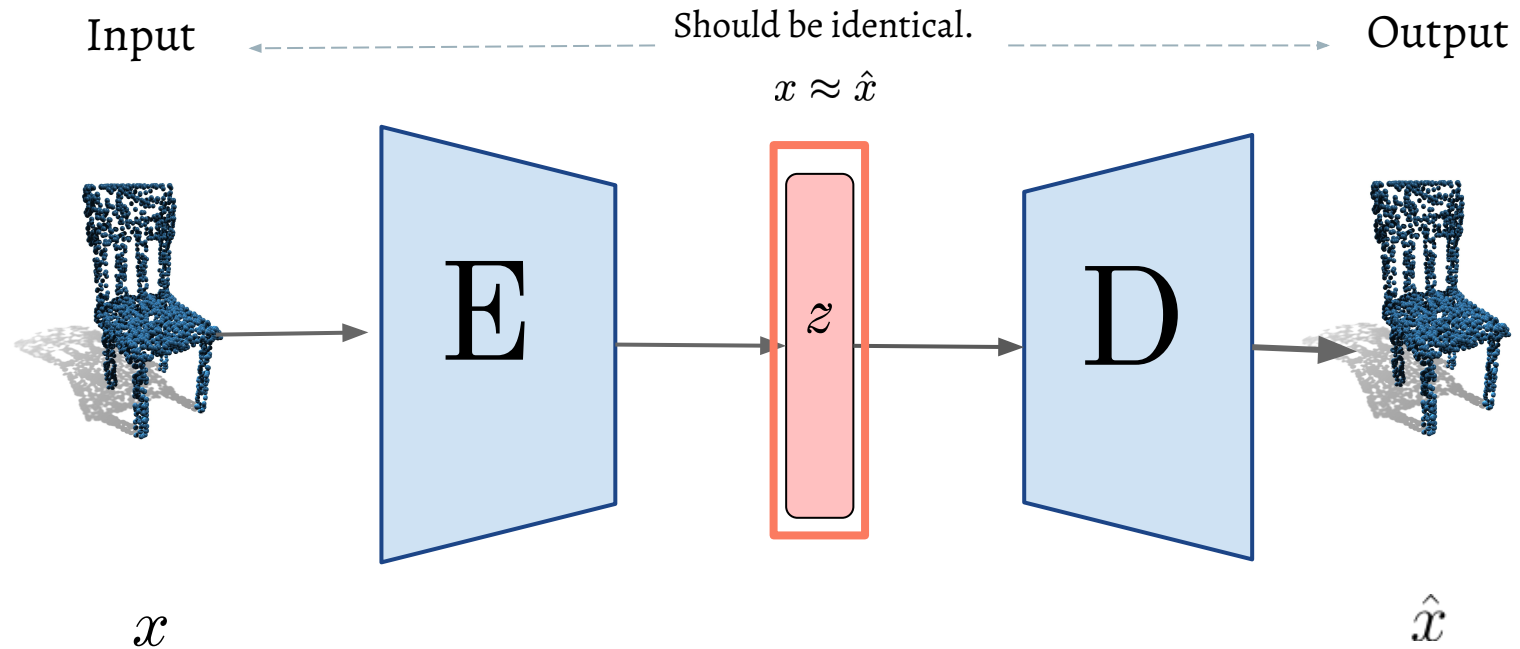
Presentation fully based on: presentation of Wojciech Stokowiec, Maciej Zięba, Maciej Zamorski, Tooploox, PLinML

Autoencoder (AE)



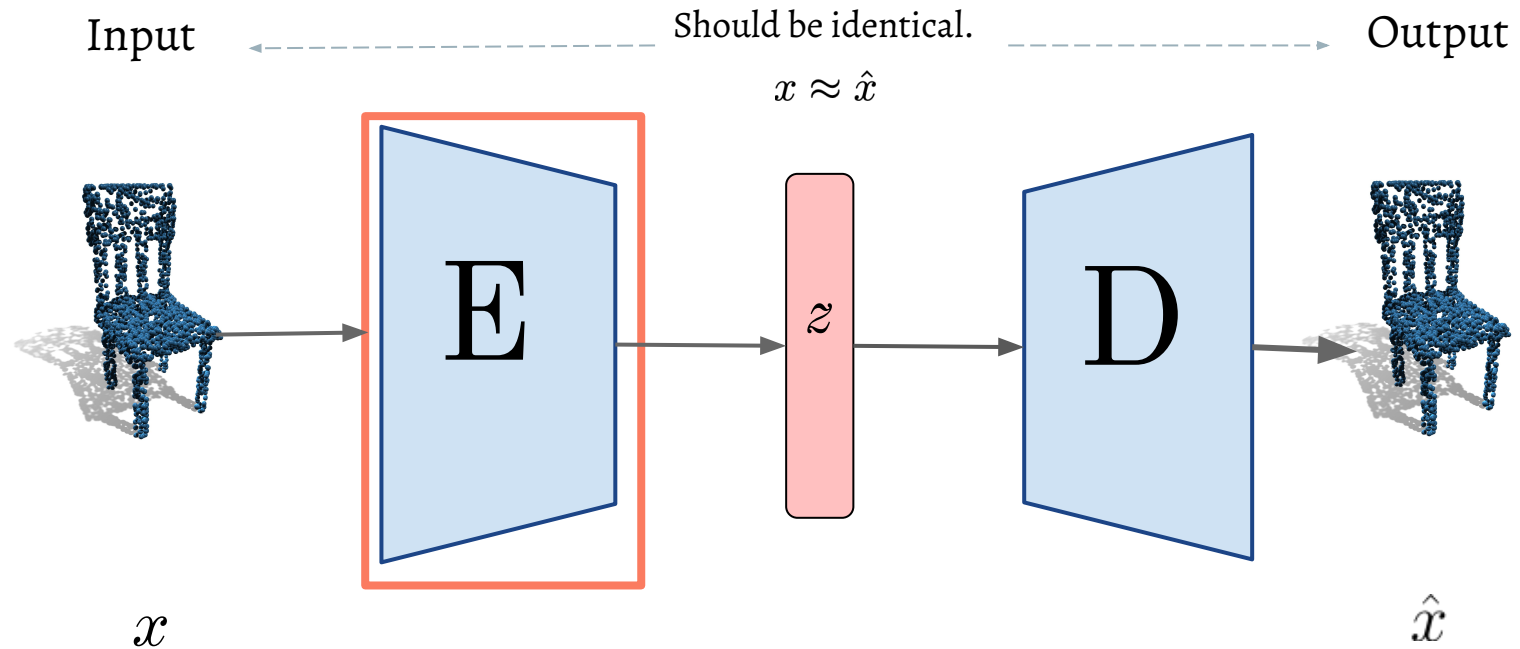
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

Autoencoder (AE)



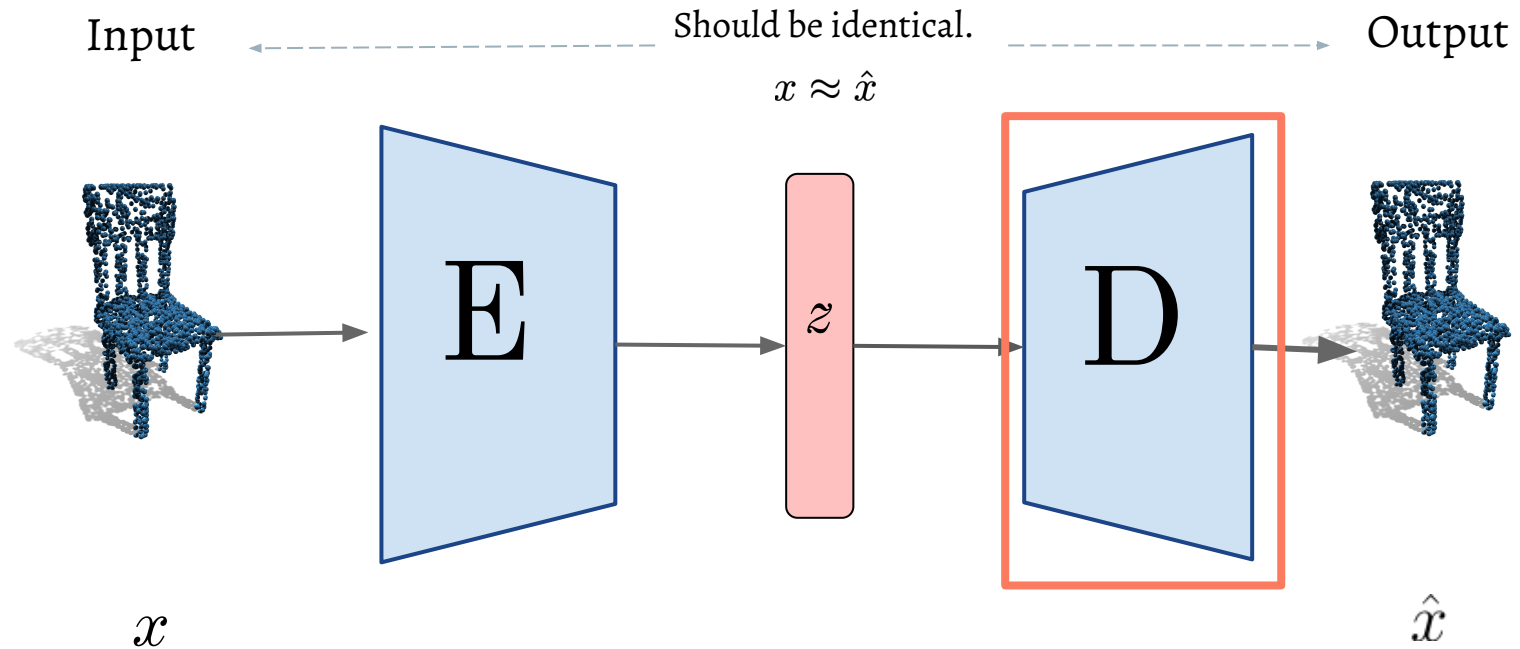
We want feature space to capture
meaningful factors of variation in
data

Autoencoder (AE)



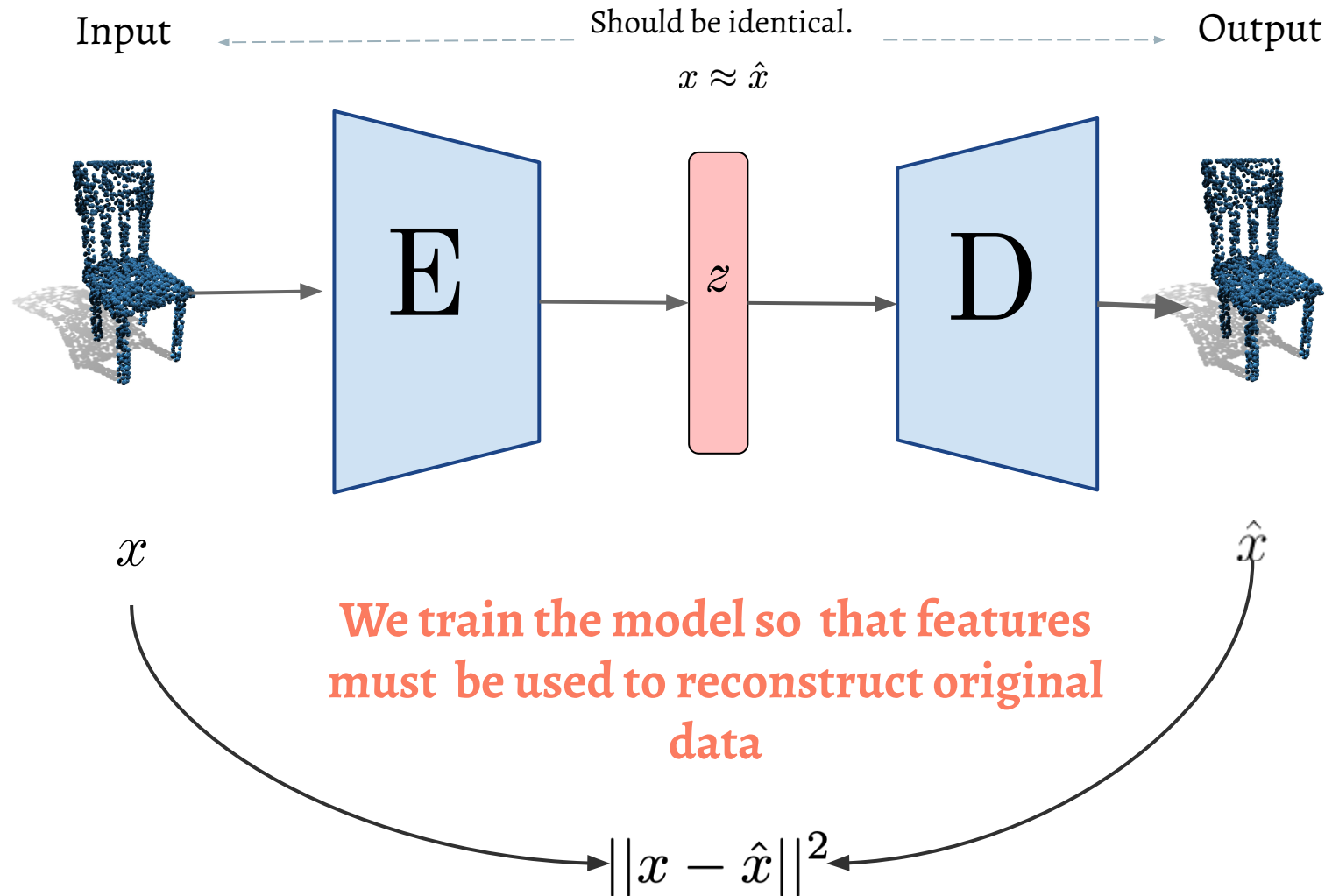
We have encoding network
(E) that encodes input
example in feature space

Autoencoder (AE)

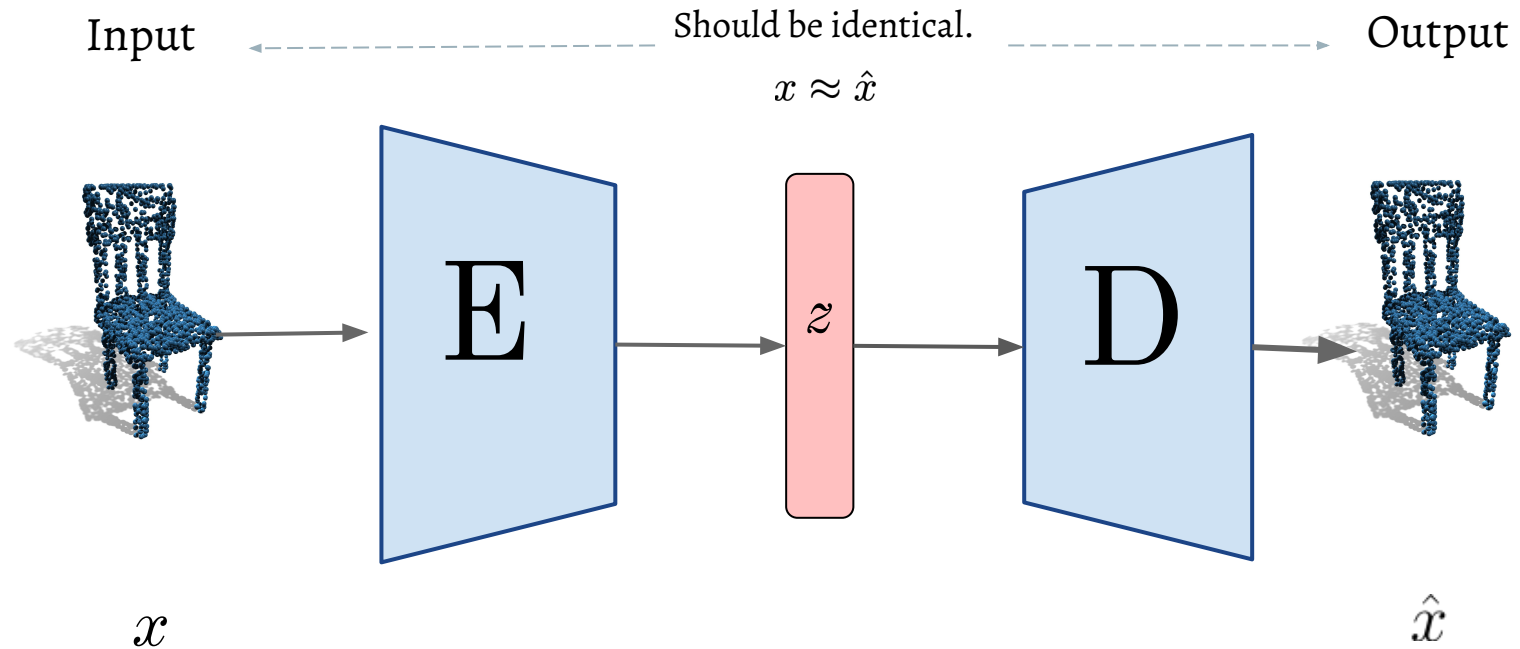


We have decoding network
(G) that restores examples
from features space

Autoencoder (AE)



Autoencoder (AE)

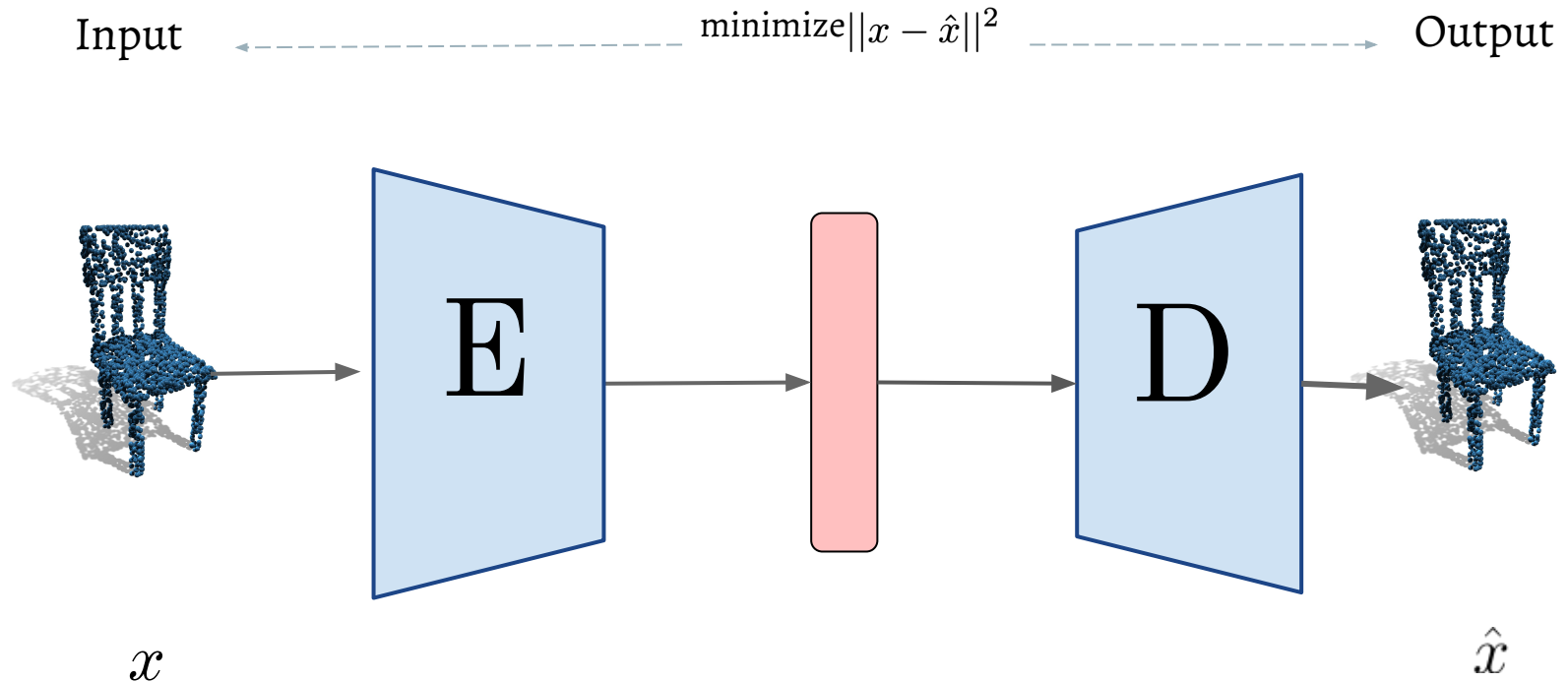


Fun fact: if linear activations are used, then optimal weights for an autoencoder can be obtained with PCA.

Variational Autoencoders

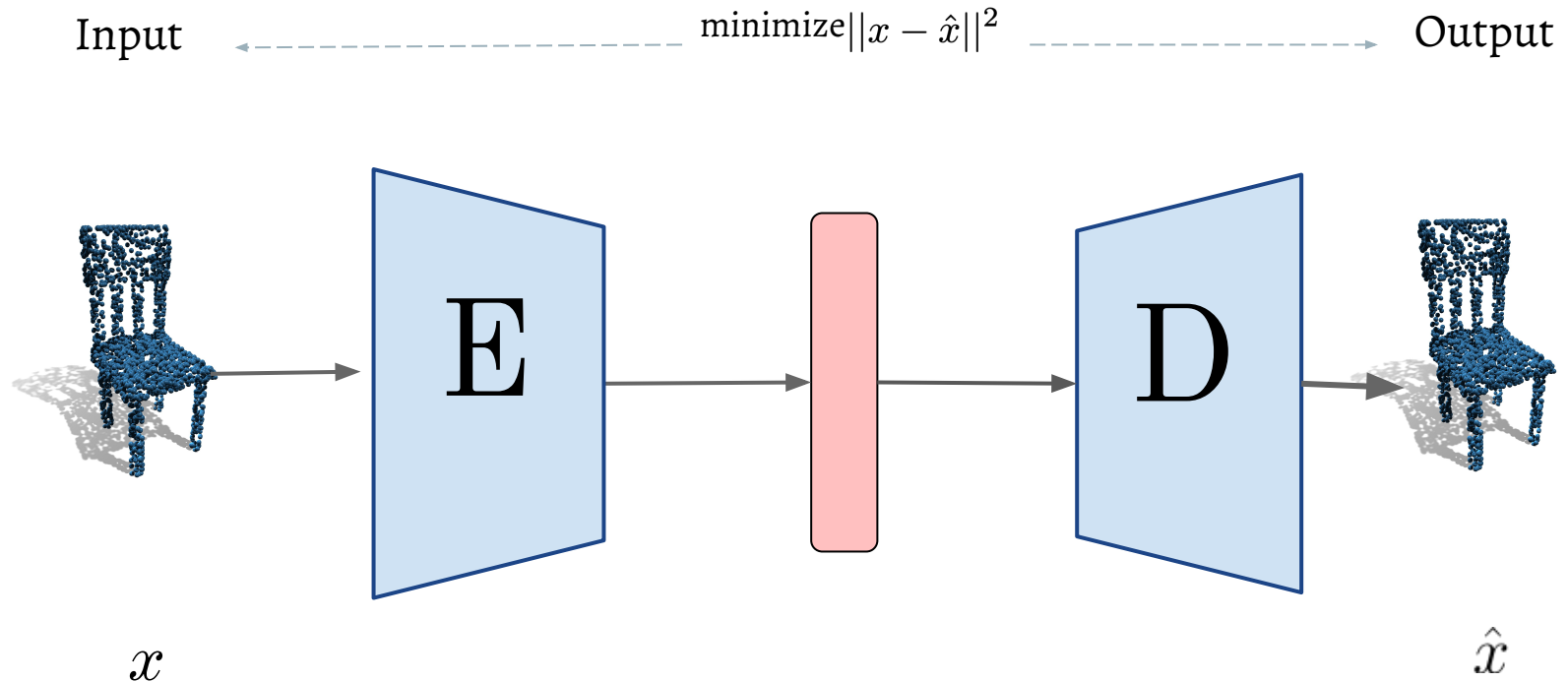
(in pictures)

Variational Autoencoder (VAE)



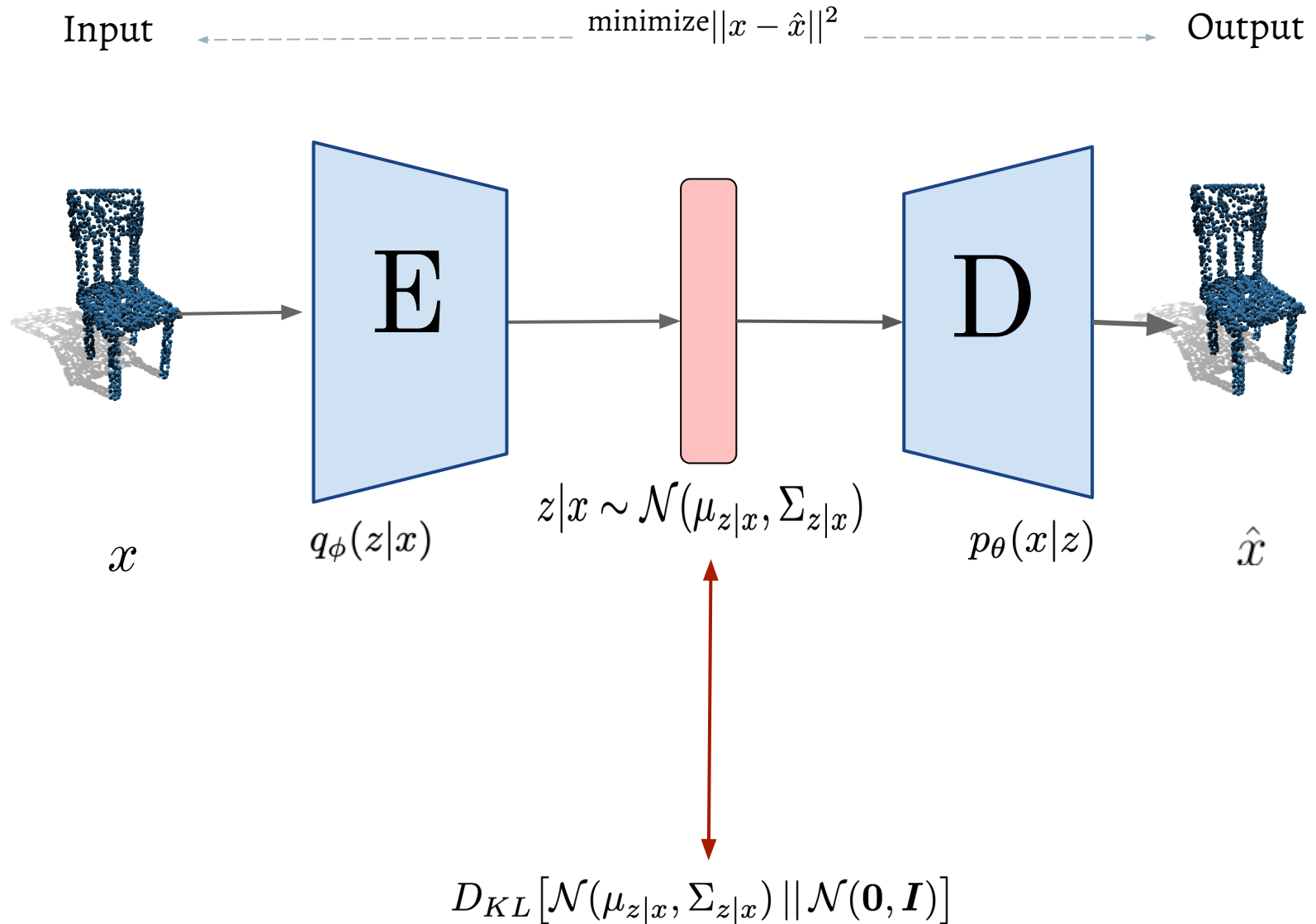
Start with a simple autoencoder ...

Variational Autoencoder (VAE)

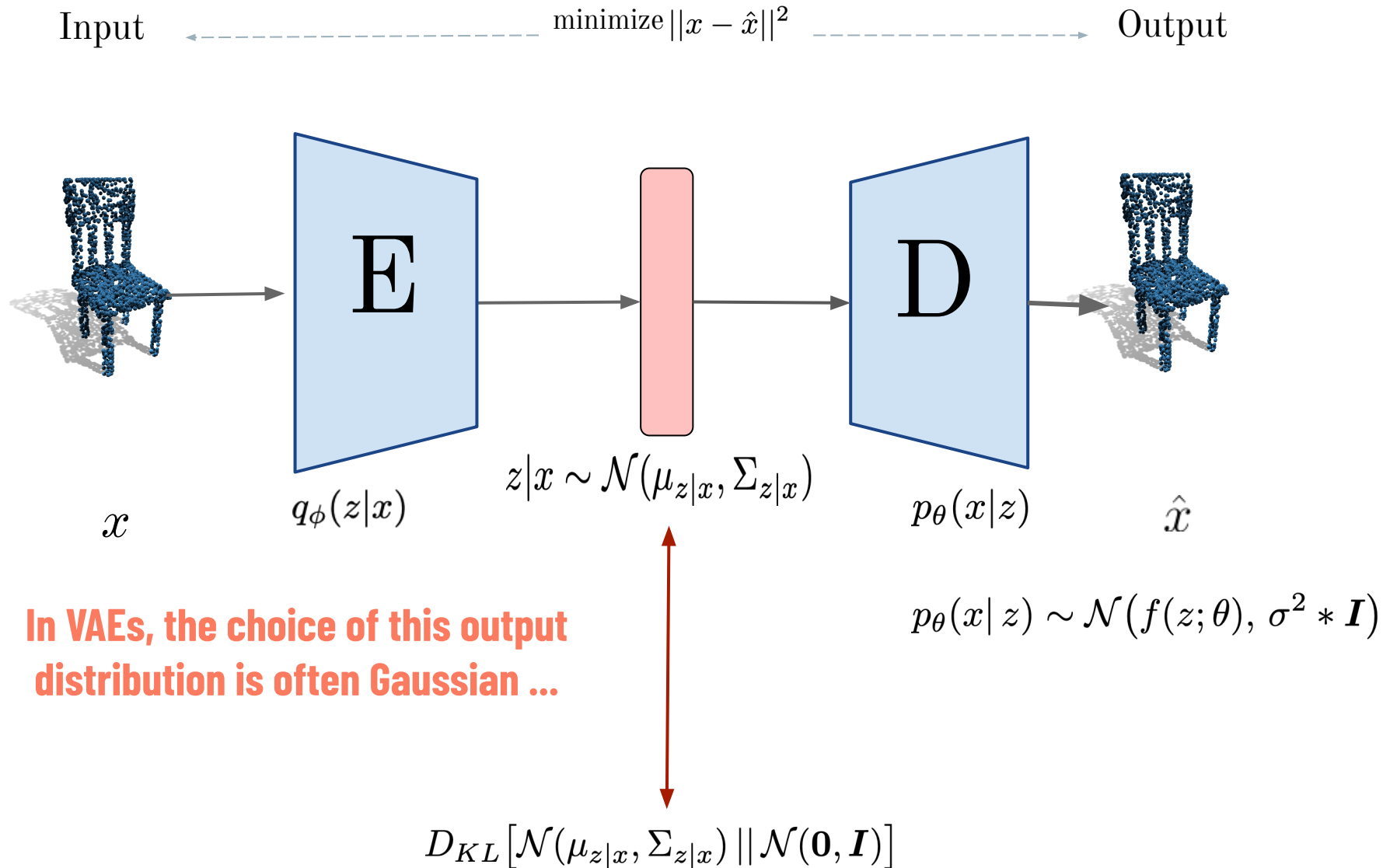


... add a probabilistic spin!

Variational Autoencoder (VAE)

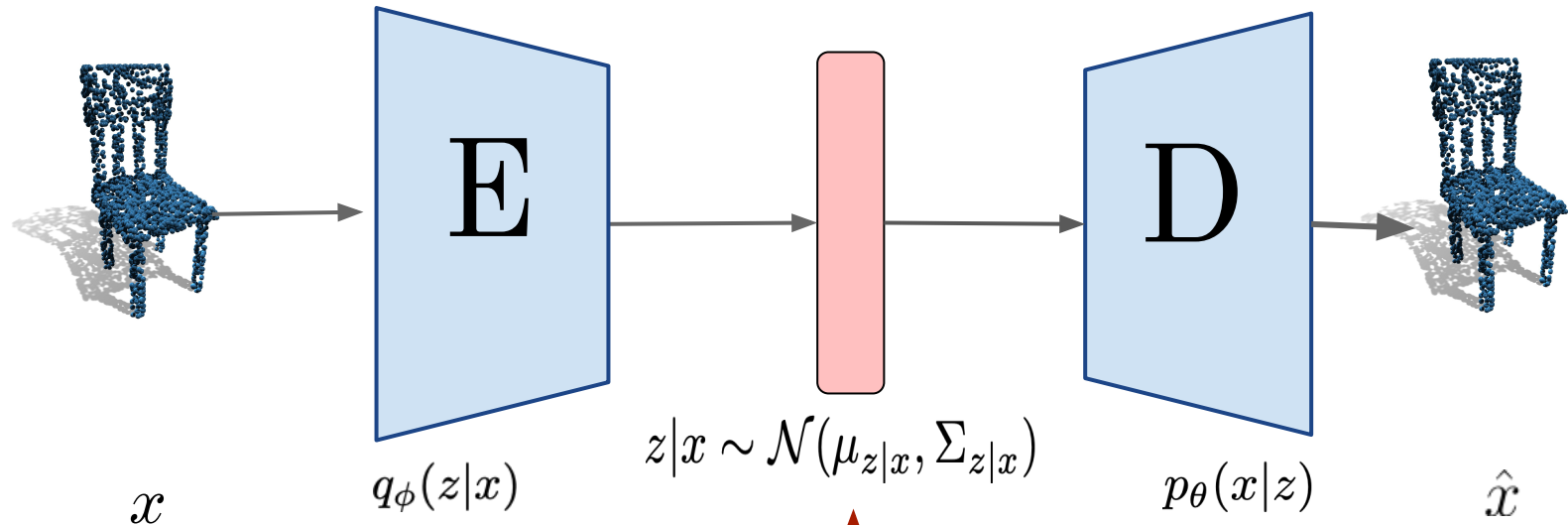


Variational Autoencoder (VAE)



Variational Autoencoder (VAE)

Input $\xleftarrow{\text{minimize } ||x - \hat{x}||^2}$ Output

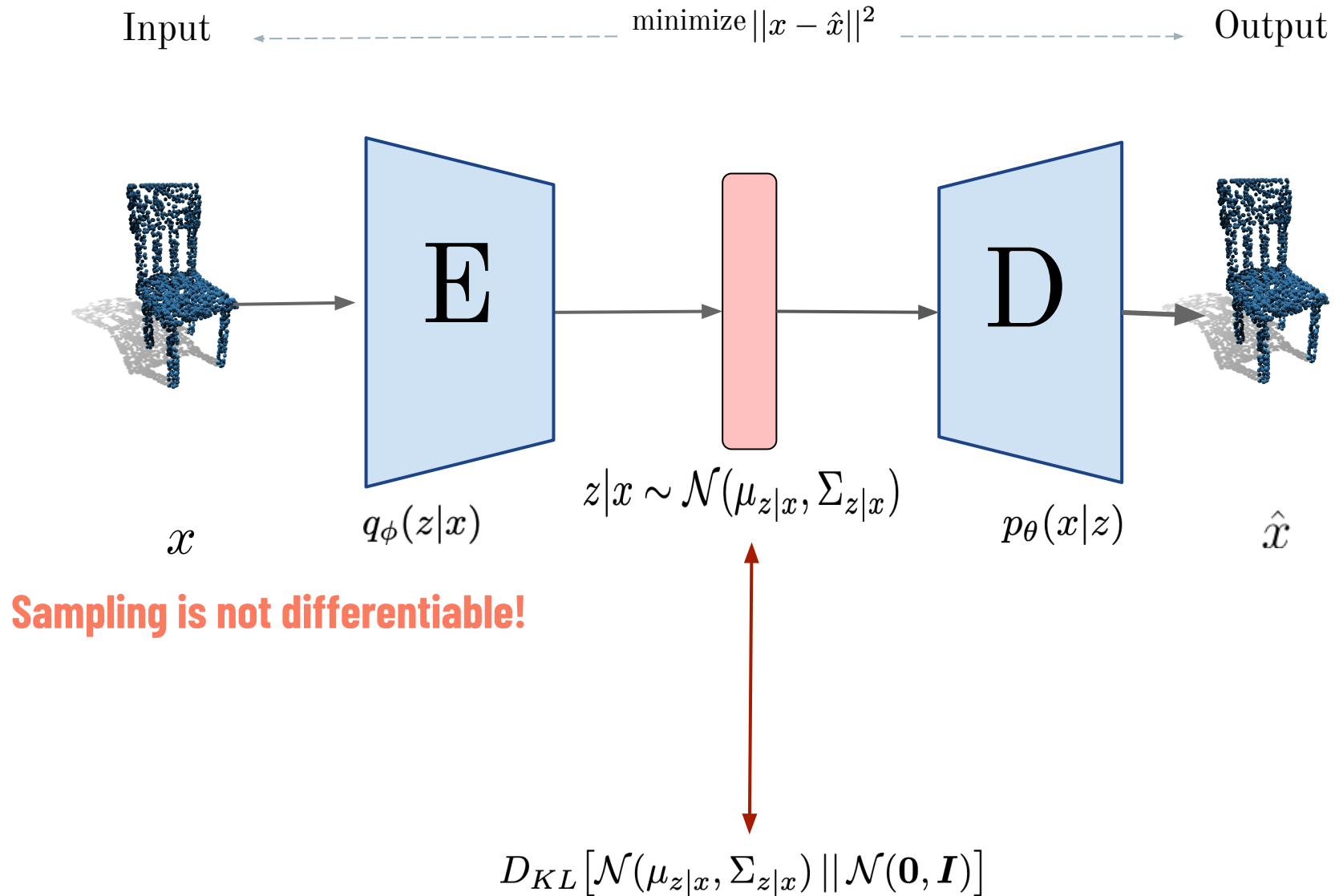


... with mean $f(z; \theta)$ and covariance equal to the identity matrix times some scalar, usually set to 1

$$p_\theta(x|z) \sim \mathcal{N}(f(z; \theta), \sigma^2 * \mathbf{I})$$

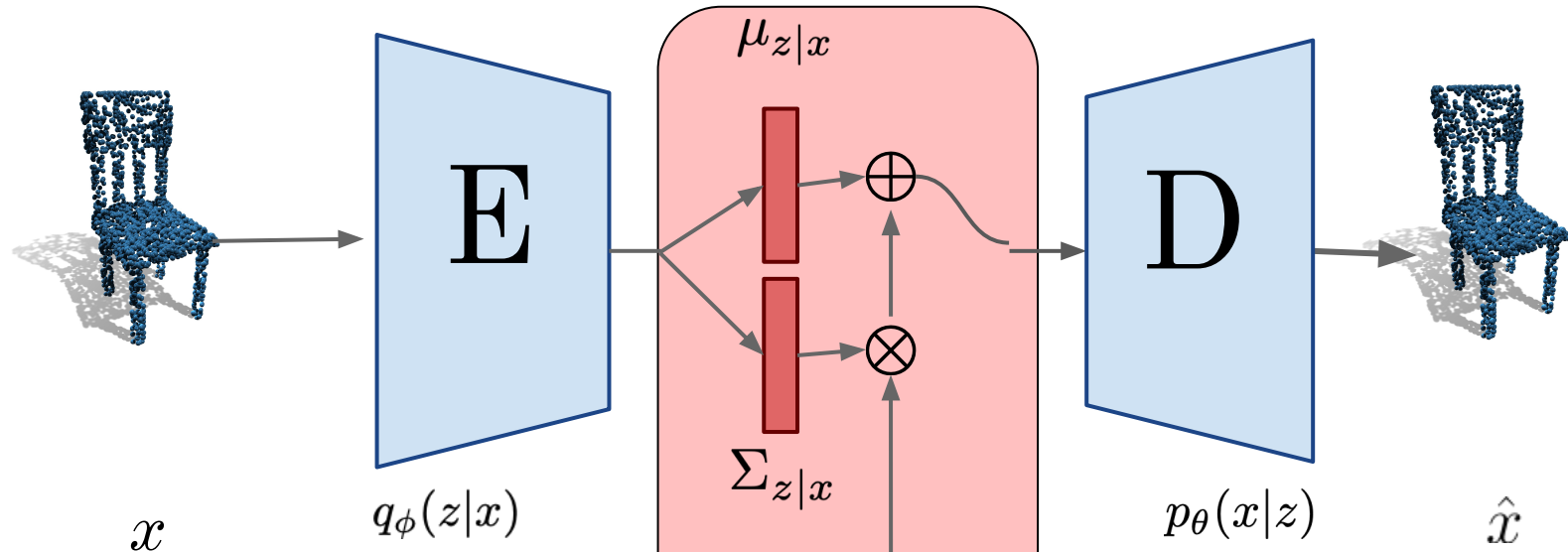
$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(\mathbf{0}, \mathbf{I})]$$

Variational Autoencoder (VAE)



Variational Autoencoder (VAE)

Input $\xleftarrow{\text{minimize } ||x - \hat{x}||^2}$ Output



To make the network end-to-end differentiable we're using re-parametrization trick

$$z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$$
$$\updownarrow$$
$$D_{KL} [\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(\mathbf{0}, \mathbf{I})]$$

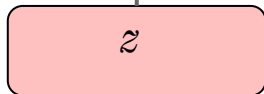
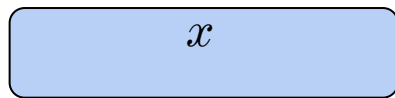
Variational Autoencoders

(in equations)

Variational Autoencoder (VAE)

We assume that training data x is generated from a latent (unobserved) random variable Z with samples denoted as z with unknown distribution.

Sample from true conditional



Sample from
true prior

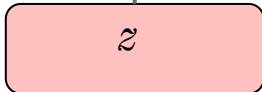
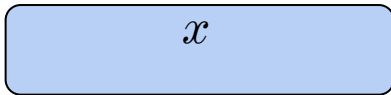
How to represent this model, to make the tractable?

- Chose prior to be simple and computationally inexpensive (gaussian).
- Approximate the conditional using neural network

Variational Autoencoder (VAE)

We assume that training data x is generated from a latent (unobserved) random variable Z with samples denoted as z with unknown distribution.

Sample from true conditional



Sample from
true prior

How to learn the parameters of the approximation of the true conditional?

- Maximize the likelihood of the training data!

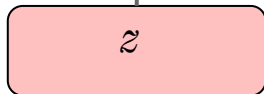
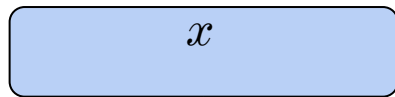
$$p(x) = \int p_{\theta}(x|z)p(z)dz$$

Would have to integrate over all the possible value of prior!

Variational Autoencoder (VAE)

We assume that training data x is generated from a latent (unobserved) random variable Z with samples denoted as z with unknown distribution.

Sample from true conditional



Sample from
true prior

How to learn the parameters of the approximation of the true conditional?

- Maximize the likelihood of the training data!

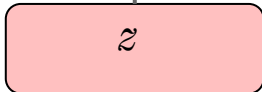
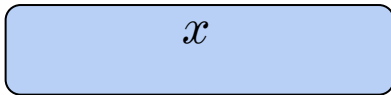
$$p(x) = \int p_{\theta}(x|z)p(z)dz$$

The key idea of VAE is to attempt to sample values of Z that are likely to have produced data and compute conditional from those

Variational Autoencoder (VAE)

We assume that training data x is generated from a latent (unobserved) random variable Z with samples denoted as z with unknown distribution.

Sample from true conditional



Sample from
true prior

How to learn the parameters of the approximation of the true conditional?

- Maximize the likelihood of the training data!

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

Add another network $q_{\phi}(z|x)$, the inference network to approximate the true posterior $p_{\theta}(z|x)$

Variational Autoencoder (VAE) - training objective

$$\begin{aligned}\log p_{\theta}(x) &= \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x) \right] \quad (p_{\theta}(x) \text{ is independent of } z) \\&= \mathbb{E}_z \left[\log \frac{p_{\theta}(x|z)p(z)}{p(z|x)} \right] \quad (\text{Bayes' Rule}) \\&= \mathbb{E}_z \left[\log \frac{p_{\theta}(x|z)p(z)}{p(z|x)} \frac{q_{\phi}(z|x)}{q_{\phi}(z|x)} \right] \quad (\text{Multiply by 1}) \\&= \mathbb{E}_z \left[\log p_{\theta}(x|z) \right] - \mathbb{E}_z \left[\log \frac{q_{\phi}(z|x)}{p(z)} \right] + \mathbb{E}_z \left[\log \frac{q_{\phi}(z|x)}{p(z|x)} \right] \quad (\text{Logarithms}) \\&= \underbrace{\mathbb{E}_z \left[\log p_{\theta}(x|z) \right] - D_{KL}(q_{\phi}(z|x) || p(z))}_{\mathcal{L}(x, \theta, \phi)} + \underbrace{D_{KL}(q_{\phi}(z|x) || p(z|x))}_{\geq 0}\end{aligned}$$

Evidence Lower Bound (**ELBO**)

$$\log p_{\theta}(x) \geq \mathcal{L}(x, \theta, \phi)$$

MLE training objective

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x_i, \theta, \phi)$$

Variational Autoencoder (VAE) - training objective

$$\log p_{\theta}(x) = \underbrace{\mathbb{E}_z [\log p_{\theta}(x | z)] - D_{KL}(q_{\phi}(z | x) || p(z))}_{\mathcal{L}(x, \theta, \phi)} + \underbrace{D_{KL}(q_{\phi}(z | x) || p(z|x))}_{\geq 0}$$

"Reconstruction error"

We want to maximize this

Approximate posterior regularization

Variational Autoencoder (VAE) - training objective

$$\log p_{\theta}(x) = \underbrace{\mathbb{E}_z [\log p_{\theta}(x | z)] - D_{KL}(q_{\phi}(z | x) || p(z))}_{\mathcal{L}(x, \theta, \phi)} + \underbrace{D_{KL}(q_{\phi}(z | x) || p(z|x))}_{\geq 0}$$

The KL divergence of the approximation and true posterior distribution, is greater or equal to zero. We can't optimize directly.

$$\log p_{\theta}(x) \geq \mathcal{L}(x_i, \theta, \phi)$$

$$\log p_{\theta}(x) \geq \mathbb{E}_z [\log p_{\theta}(x | z)] - D_{KL}(q_{\phi}(z | x) || p(z))$$

Variational Autoencoder (VAE) - reconstruction term

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x | z) \right]$$

This called the reconstruction term

Let see why training loss is $\|x - \hat{x}\|^2$ while the objective talks about expectations and log probabilities

Variational Autoencoder (VAE) - reconstruction term

we decided to model the conditional distribution as multivariate gaussian

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log p_{\theta}(x|z) \right] p_{\theta}(x|z) \sim \mathcal{N}(f(z; \theta), \sigma^2 * \mathbf{I})$$

Recall that

$$p(x|\mu, \Sigma) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

$$\ln p(x|\mu, \Sigma) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + C$$

Variational Autoencoder (VAE) - reconstruction term

we decided to model the conditional distribution as multivariate gaussian

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] p_{\theta}(x|z) \sim \mathcal{N}(f(z; \theta), \sigma^2 * \mathbf{I})$$

$$\ln p(x|\mu, \Sigma) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) + C$$

$$\ln p(x|\mu, \sigma^2 I) = -\frac{1}{2} \ln |\sigma^2 I| - \frac{1}{2} (x - \mu)^T (\sigma^2 I)^{-1} (x - \mu) + C$$

Only L-2 norm is
dependent on model
parameters!

$$= \dots$$

$$= -\frac{1}{2} \ln \sigma^2 - \frac{1}{2} \|x - \mu\|^2 / \sigma^2 + C$$

$$\ln p_{\theta}(x|f(z; \theta), I) = -\frac{1}{2} \|x - f(z; \theta)\|^2 + C$$

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}(q_{\phi}(z | x) || p(z))$$

How can we calculate KL divergence between the approximate posterior and prior?

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I)]$$

Both approximate posterior and prior are assumed to be multivariate gaussian, we can do this analytically!

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I)] =$$
$$\frac{1}{2} \left(\text{tr}(\Sigma_{z|x}) + \mu_{z|x}^T \mu_{z|x} - k - \log \det(\Sigma_{z|x}) \right)$$

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma_{z|x}) + \mu_{z|x}^T \mu_{z|x} - k - \log \det(\Sigma_{z|x}) \right)$$

Easy, just sum!

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma_{z|x}) + \boxed{\mu_{z|x}^T \mu_{z|x}} - k - \log \det(\Sigma_{z|x}) \right)$$

Easy, simple
dot product!

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) \parallel \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma_{z|x}) + \mu_{z|x}^T \mu_{z|x} - \boxed{k} - \log \det(\Sigma_{z|x}) \right)$$

Easy, scalar.
Equal to dimensionality of the distribution

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) || \mathcal{N}(0, I)] = \frac{1}{2} \left(\text{tr}(\Sigma_{z|x}) + \mu_{z|x}^T \mu_{z|x} - k - \log \det(\Sigma_{z|x}) \right)$$

Easy, as we assumed that the covariance matrix of approximate posterior is diagonal.

Variational Autoencoder (VAE) - Kullback-Leibler divergence

$$D_{KL}[\mathcal{N}(\mu_{z|x}, \Sigma_{z|x}) \parallel \mathcal{N}(0, I)] =$$
$$\frac{1}{2} \sum_{i=1}^k (\sigma_i^2 + \mu_i^2 - \log \sigma_i^2 - 1)$$

Variational Autoencoder (VAE) - objective summary

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} [q_{\phi}(z|x) || p(z)]$$

We train the model to maximize this expression consisting of two terms:

- Reconstruction term
- Regularization term on the approximate posterior

Variational Autoencoder (VAE) - objective summary

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} [q_{\phi}(z|x) || p(z)]$$

Is maximized by minimizing the
reconstruction error

Variational Autoencoder (VAE) - objective summary

$$\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} [q_{\phi}(z|x) || p(z)]$$

Has a simple analytical formula, that can be easily minimized using gradient methods

Variational Autoencoder (VAE) - objective summary

$$\underbrace{\mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)]}_{-\frac{1}{2} ||x - \hat{x}||^2 + C} - \underbrace{D_{KL} [q_{\phi}(z|x) || p(z)]}_{\frac{1}{2} \sum_{i=1}^k (\sigma_i^2 + \mu_i^2 - \log \sigma_i^2 - 1)}$$

ELBO in theory vs ELBO in “almost code”

Normalizing flows

Tomasz Kajdanowicz, Piotr Bielek, Maciej Falkiewicz,
Kacper Kania, Piotr Zieliński

Rereferences

1. Tabak, E. G., & Turner, C. V. (2013). A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics*, 66(2), 145-164.
2. **Dinh, L., Krueger, D., & Bengio, Y. (2014). Nice: Non-linear independent components estimation. arXiv preprint arXiv:1410.8516.**
3. **Rezende, D. J., & Mohamed, S. (2015). Variational inference with normalizing flows. arXiv preprint arXiv:1505.05770.**
4. **Dinh, Laurent, Jascha Sohl-Dickstein, and Samy Bengio. "Density estimation using real nvp." arXiv preprint arXiv:1605.08803 (2016).**

Further reading:

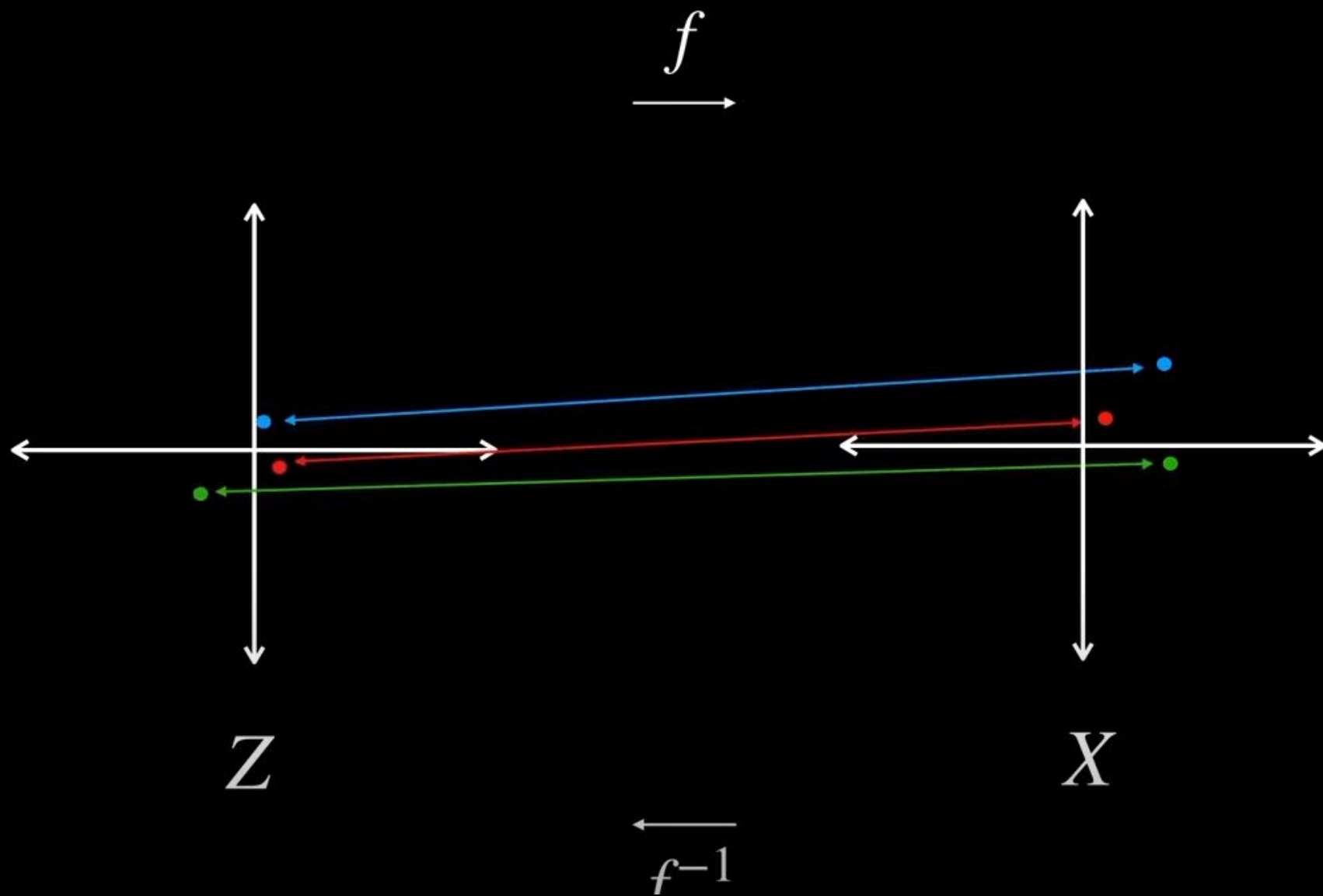
Adam Kosiosek blog post - http://akosiosek.github.io/ml/2018/04/03/norm_flows.html

Presentation fully based on: presentation of Ari Seff, Princeton University,
<https://www.youtube.com/watch?v=i7LjDvsLWCg>

$$z \sim p_{\theta}(z) = \mathcal{N}(z; 0, I)$$

$$x = f_{\theta}(z) = f_K \circ \dots \circ f_2 \circ f_1(z)$$

each f_i is invertible (bijective)



$$z \sim p_{\theta}(z) = \mathcal{N}(z; 0, I)$$

$$x = f_{\theta}(z) = f_K \circ \dots \circ f_2 \circ f_1(z)$$

$$p_{\theta}(x) \stackrel{?}{=} p_{\theta}(f_{\theta}^{-1}(x))$$

Change of variables formula

$f: Z \rightarrow X$, f is invertible

$p_\theta(z)$ defined over $z \in Z$

Change of variables formula:

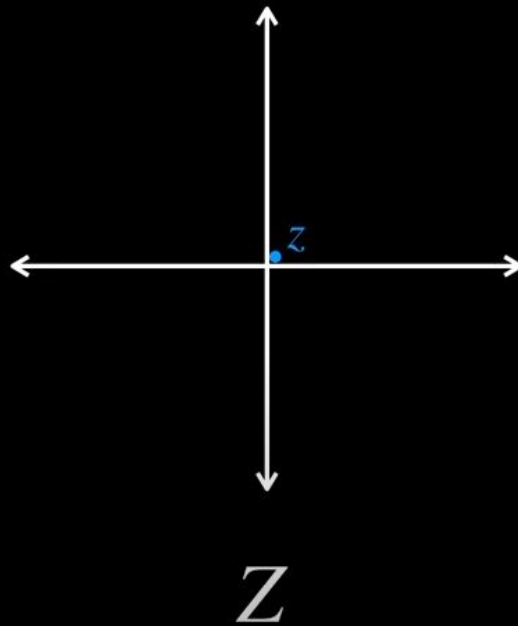
$$p_\theta(x) = p_\theta(f^{-1}(x)) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

$$p_{\theta}(x) = p_{\theta}(f^{-1}(x)) \left| \det \left(\frac{\partial f^{-1}(x)}{\partial x} \right) \right|$$

$$= p_{\theta}(z) \left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$

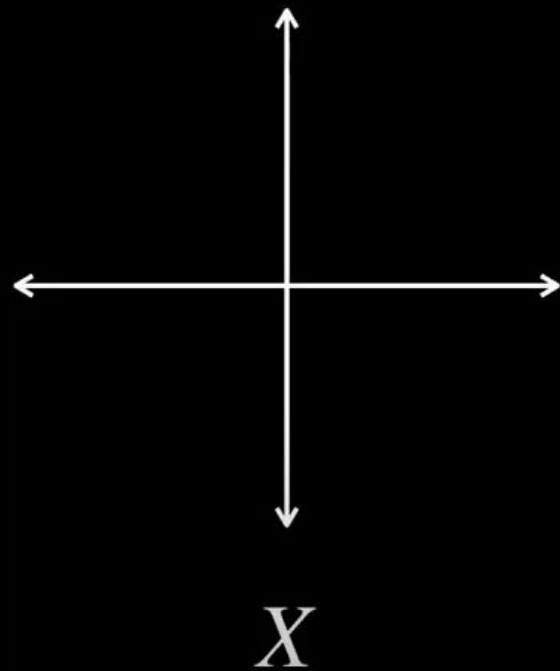
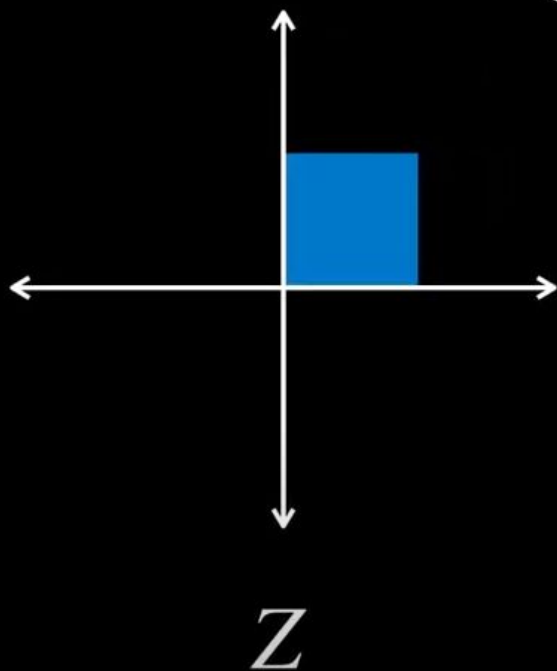
$$p_{\theta}(x) = p_{\theta}(z) \left| \det \left(\frac{\partial z}{\partial x} \right) \right|$$

Example



Example

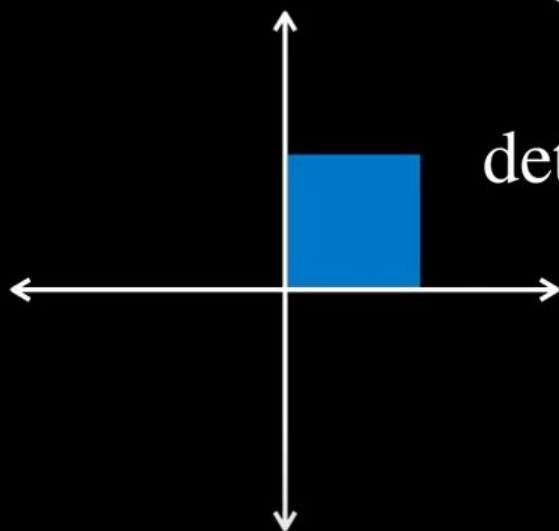
$$J = \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \frac{\partial x_1}{\partial z_2} \\ \frac{\partial x_2}{\partial z_1} & \frac{\partial x_2}{\partial z_2} \end{bmatrix}$$



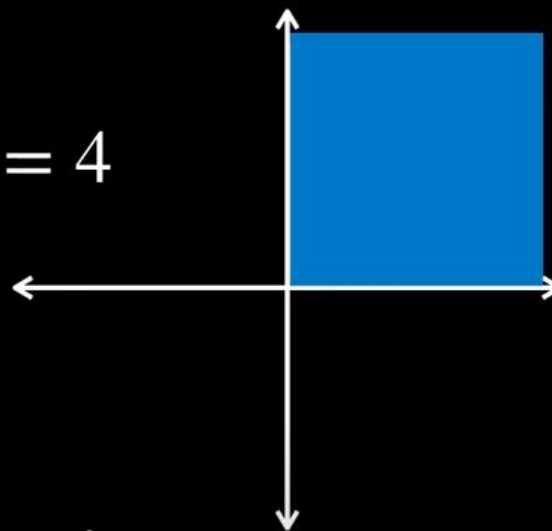
Example

$$J = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

$$\det(J) = 2 \cdot 2 - 0 \cdot 0 = 4$$



z

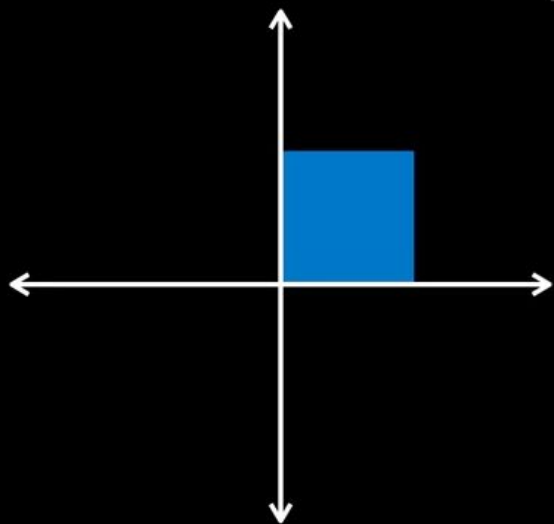


x

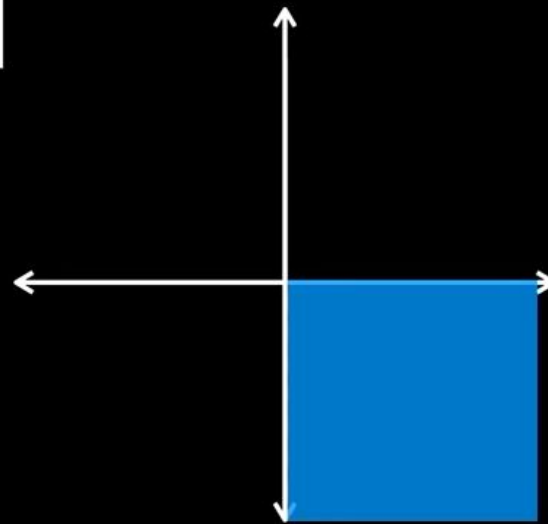
$$p(x) = p(z) \left| \frac{1}{\det\left(\frac{\partial x}{\partial z}\right)} \right|$$

Example

$$J = \begin{bmatrix} 2 & 0 \\ 0 & -2 \end{bmatrix}$$



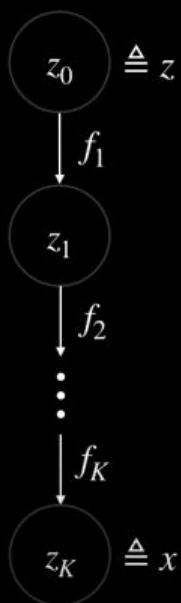
z



x

$$p(x) = p(z) \left| \frac{1}{\det\left(\frac{\partial x}{\partial z}\right)} \right|$$

Normalizing flow



$$z \sim p_\theta(z)$$

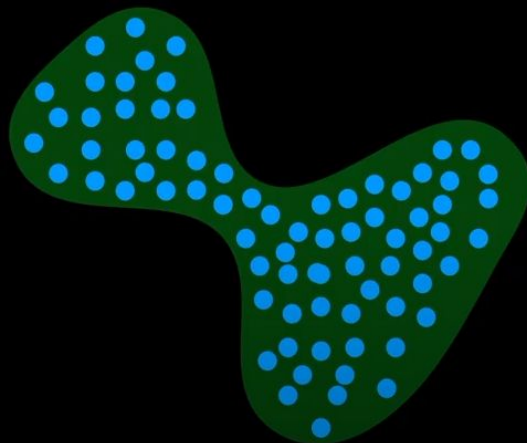
“normalizing flow”

$$x = f_\theta(z) = f_K \circ \dots \circ f_2 \circ f_1(z)$$

$$p_\theta(x) = p_\theta(z) \prod_1^K \left| \det \left(\frac{\partial f_i^{-1}}{\partial z_i} \right) \right| = p_\theta(z) \left| \det \left(\frac{\partial f^{-1}}{\partial x} \right) \right|$$

Applications

$p_D(x)$



$p_\theta(x)$

$p_\theta(x | z)$ (likelihood)

z



z



z



z

$p_\theta(z)$ (prior)

observed
variables

x

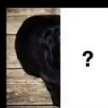


$p_\theta(x)$

x



x



p_θ



vs.

p_θ



How to use normalizing flows?

$$z \sim p_{\theta}(z)$$

$$x = f_{\theta}(z) = f_K \circ \dots \circ f_2 \circ f_1(z)$$

$$p_{\theta}(x) = p_{\theta}(z) \left| \det \left(\frac{\partial f^{-1}}{\partial x} \right) \right|$$

How to use normalizing flows?

$$p_{\theta}(x) = p_{\theta}(z) \left| \det \left(\frac{\partial f^{-1}}{\partial x} \right) \right|$$

$$\log p_{\theta}(x) = \log p_{\theta}(z) + \log \left| \det \left(\frac{\partial f^{-1}}{\partial x} \right) \right|$$

$$= \log p_{\theta}(z) + \sum_{i=1}^K \log \left| \det \left(\frac{\partial f_i^{-1}}{\partial z_i} \right) \right|$$

How to use normalizing flows?

$$\log p_{\theta}(x) = \log p_{\theta}(z) + \sum_{i=1}^K \log \left| \det \left(\frac{\partial f_i^{-1}}{\partial z_i} \right) \right|$$

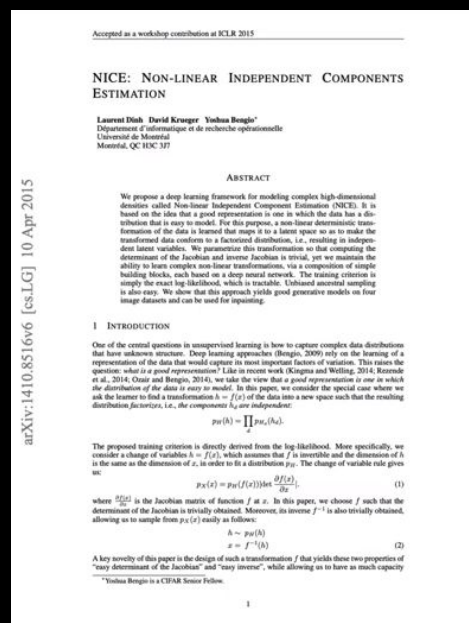
exact log-likelihood evaluation

exact posterior inference (via $z = f^{-1}(x)$)

Jacobian determinant computations are expensive

$$\log p_{\theta}(x) = \log p_{\theta}(z) + \sum_{i=1}^K \log \left| \det \left(\frac{\partial f_i^{-1}}{\partial z_i} \right) \right|$$

$$\begin{bmatrix} a_{1,1} & 0 & \dots & \dots & 0 \\ a_{2,1} & a_{2,2} & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ a_{n,1} & a_{n,2} & \dots & \dots & a_{n,n} \end{bmatrix}$$

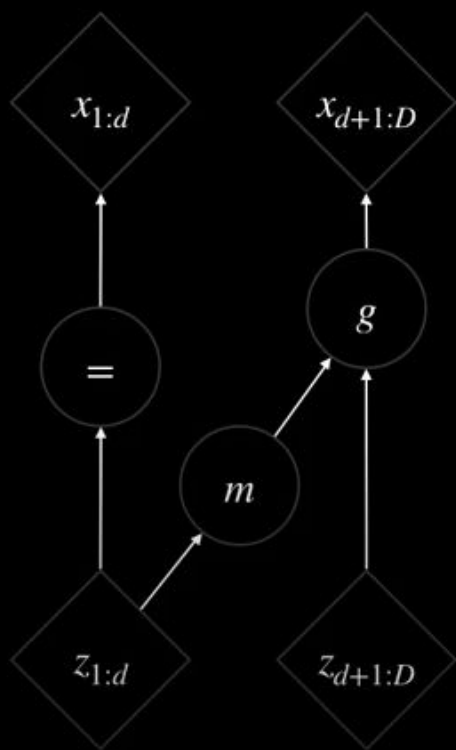


NICE - nonlinear independent components estimation

Coupling layers

$$x_{1:d} = z_{1:d}$$

$$x_{d+1:D} = g(z_{d+1:D}; m(z_{1:d}))$$



$$\frac{\partial x}{\partial z} = \begin{bmatrix} I_d & 0 \\ \frac{\partial x_{d+1:D}}{\partial z_{1:d}} & \frac{\partial x_{d+1:D}}{\partial z_{d+1:D}} \end{bmatrix}$$

Inversion of transformation

$$x_{1:d} = z_{1:d}$$

$$x_{d+1:D} = g(z_{d+1:D}; m(z_{1:d}))$$

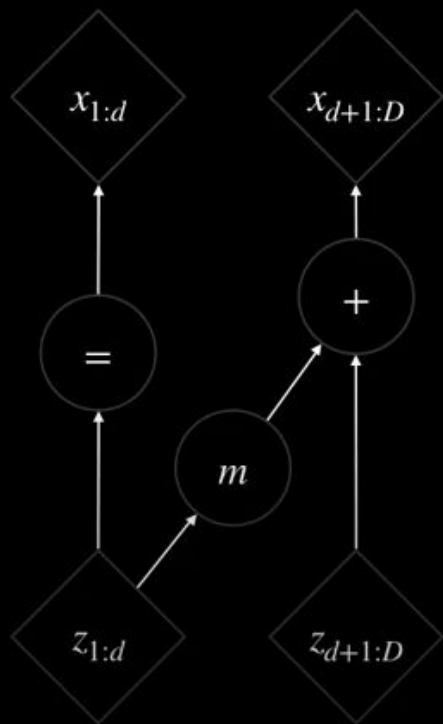


$$z_{1:d} = x_{1:d}$$

$$z_{d+1:D} = g^{-1}(x_{d+1:D}; m(x_{1:d}))$$

Additive coupling layer

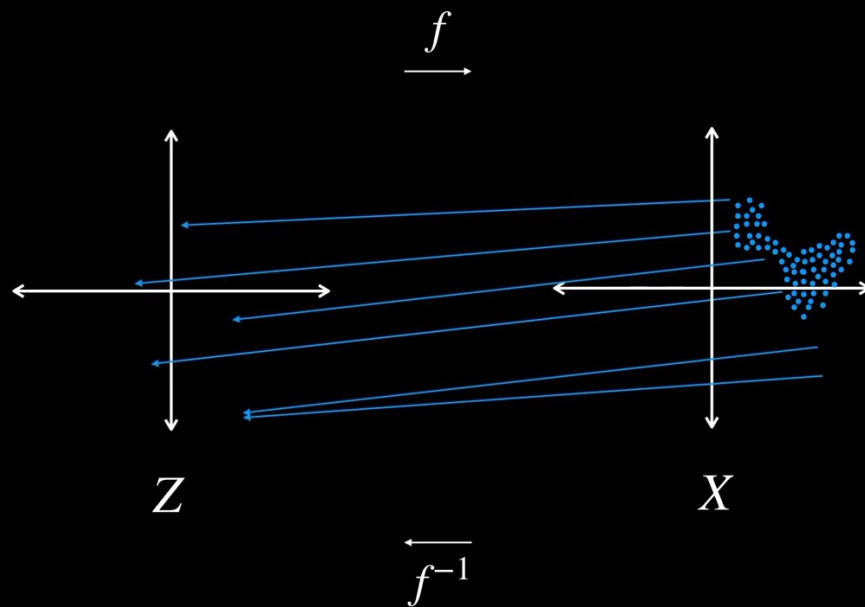
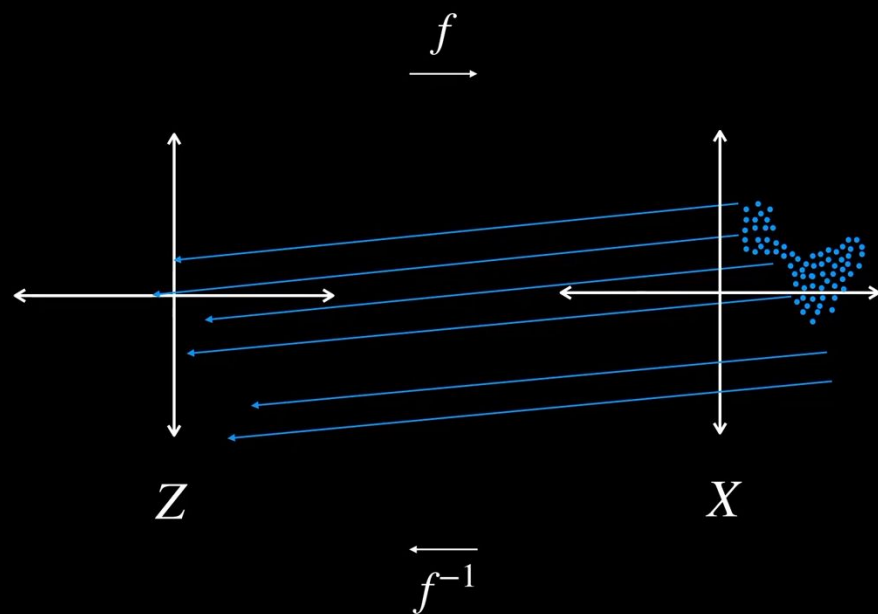
$$g(z_{d+1:D}; m(z_{1:d})) = z_{d+1:D} + m(z_{1:d})$$



$$x_{1:d} = z_{1:d}$$

$$x_{d+1:D} = z_{d+1:D} + m(z_{1:d})$$

Limitations



Scaling matrix

$$\begin{bmatrix} S_{1,1} & 0 & \dots & \dots & 0 \\ 0 & S_{2,2} & 0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & 0 & \dots & \dots & S_{D,D} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_D \end{bmatrix}$$

Scaling matrix S

$$\log p_{\theta}(x) = \sum_{i=1}^D [\log(p_{\theta}(f^{-1}(x)_i)) - \log(|S_{ii}|)]$$

DENSITY ESTIMATION USING REAL NVP

Laurent Dinh*
 Montreal Institute for Learning Algorithms
 University of Montreal
 Montreal, QC H3T 1J4

Jaehyeon Kim*
 Google Brain

Samy Bengio
 Google Brain

ABSTRACT

Unsupervised learning of probabilistic models is a central yet challenging problem in machine learning. Specifically, designing models with tractable learning, sampling, inference and evaluation is crucial to solving this task. We extend the space of such models using real-valued non-volume preserving (real NVP) transformations, a set of powerful, easily invertible, and invertible transformations, resulting in an ensemble model training algorithm with exact log-likelihood computation, exact and efficient sampling, exact and efficient inference, tractable gradients, and an interpretable latent space. We demonstrate its ability to model natural images on four datasets through sampling, log-likelihood evaluation, and latent variable manipulations.

1 Introduction

The domain of representation learning has undergone tremendous advances due to improved unsupervised learning techniques. However, unsupervised learning has the potential to leverage large pools of unlabeled data, and convert these advances to modifying. But are otherwise impractical or infeasible. One principled approach to unsupervised learning is generative probabilistic modeling. Not only do generative probabilistic models have the ability to create novel content, they also have a rich range of reconstruction related applications including inpainting [31, 46, 39], denoising [1], colorization [17], and super-resolution [5].

As data of interest are generally high-dimensional and highly structured, the challenge in this domain is building models that are powerful enough to capture its complexity yet still tractable. We address this challenge by introducing real-valued non-volume preserving (real NVP) transformations, a tractable yet expressive approach to modeling high-dimensional data.

This model can perform efficient and exact inference, sampling and log density estimation of data points. Moreover, the architecture presented in this paper enables exact and efficient reconstruction of input images from the hierarchical features extracted by this model.

2 Related work

Substantial work on probabilistic generative models has focused on training models using maximum likelihood. One class of maximum likelihood models are those described by probabilistic, undirected graphs, such as Restricted Boltzmann Machines [26] and Deep Boltzmann Machines [52]. These models are trained by taking advantage of the conditional independence property of their bipartite structure to allow efficient exact or approximate gradient estimation on latent variables. However, because of the intractability of the associated marginal distribution over latent variables, their training, evaluation, and sampling procedures necessitate the use of approximations. See *Mean Field inference* and *Markov Chain Monte Carlo*, whose convergence time for such complex models

*Work was done while author was at Google Brain.



z



x



z

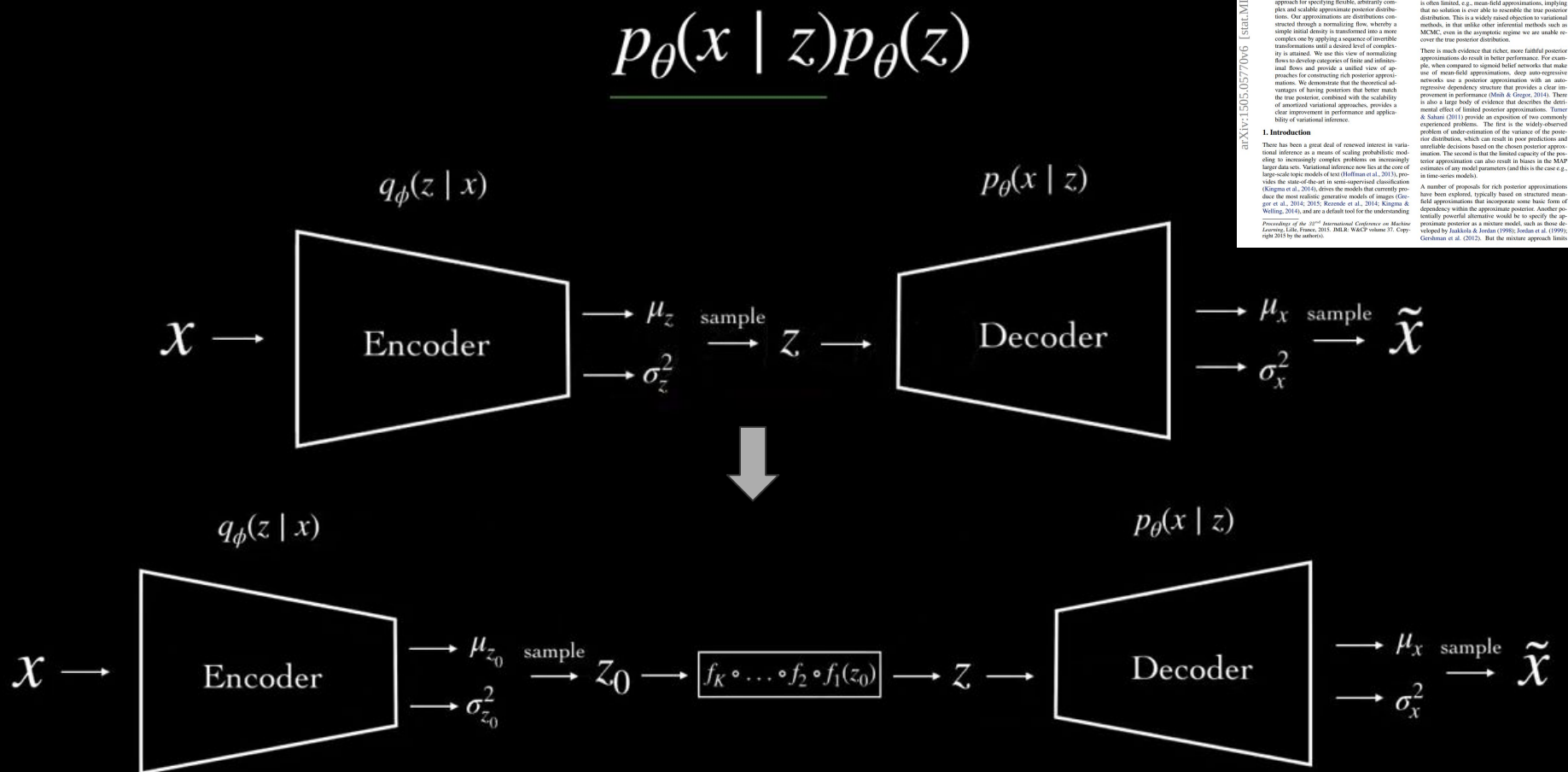


x

$$z = (z^{(1)}, \dots, z^{(L)})$$

each $z^{(i)}$ operates at a different scale

Other applications



Summary

Normalizing flows:

- handle complex probability distributions
- there are some tricks to make it computationally tractable
- only a few of the ways of using them shown in the presentation