# Introduction to Machine Learning: Lecture III

## Michael Kagan

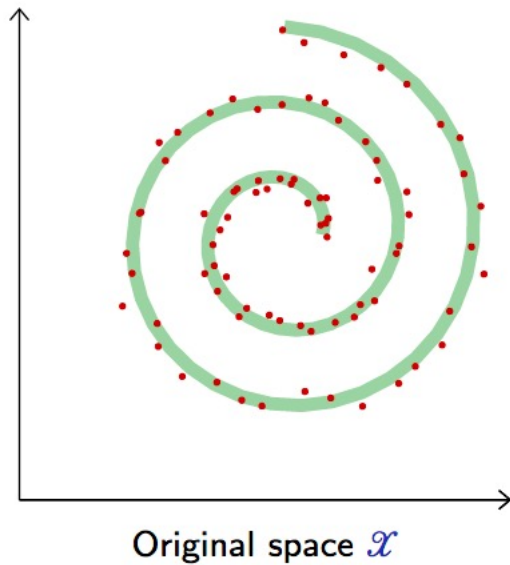### SLAC

INFN School of Statistics 2022
May 19, 2022

# The Plan
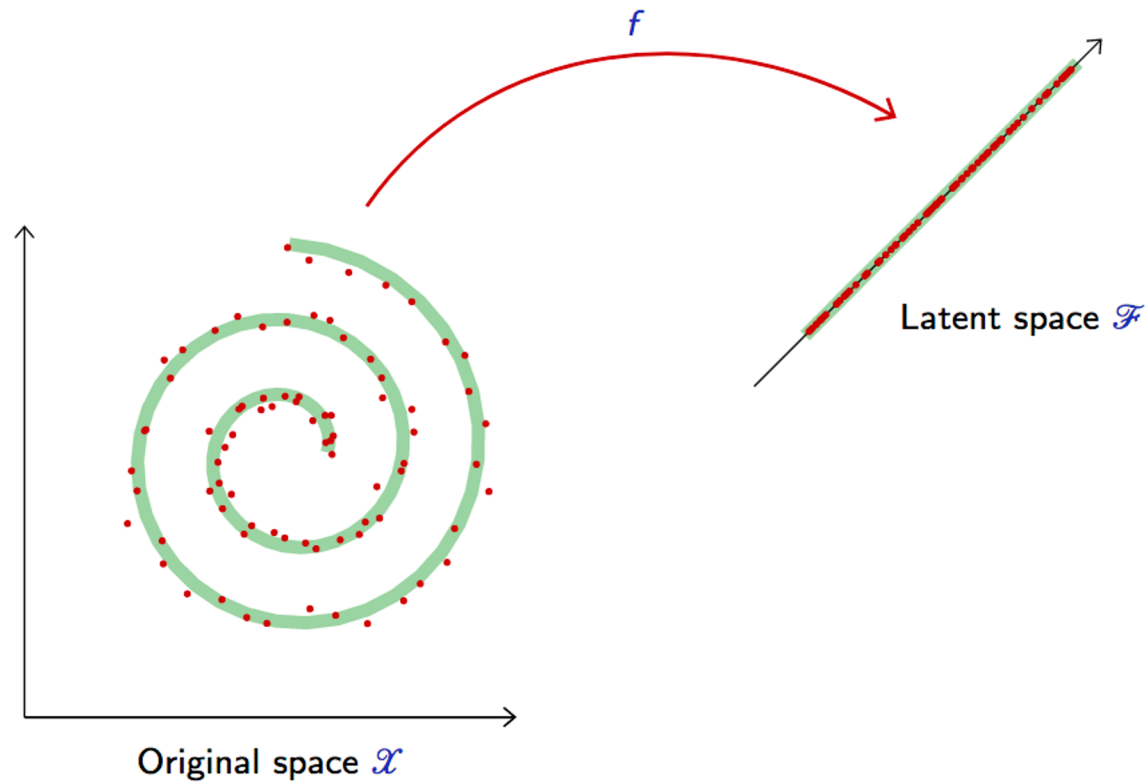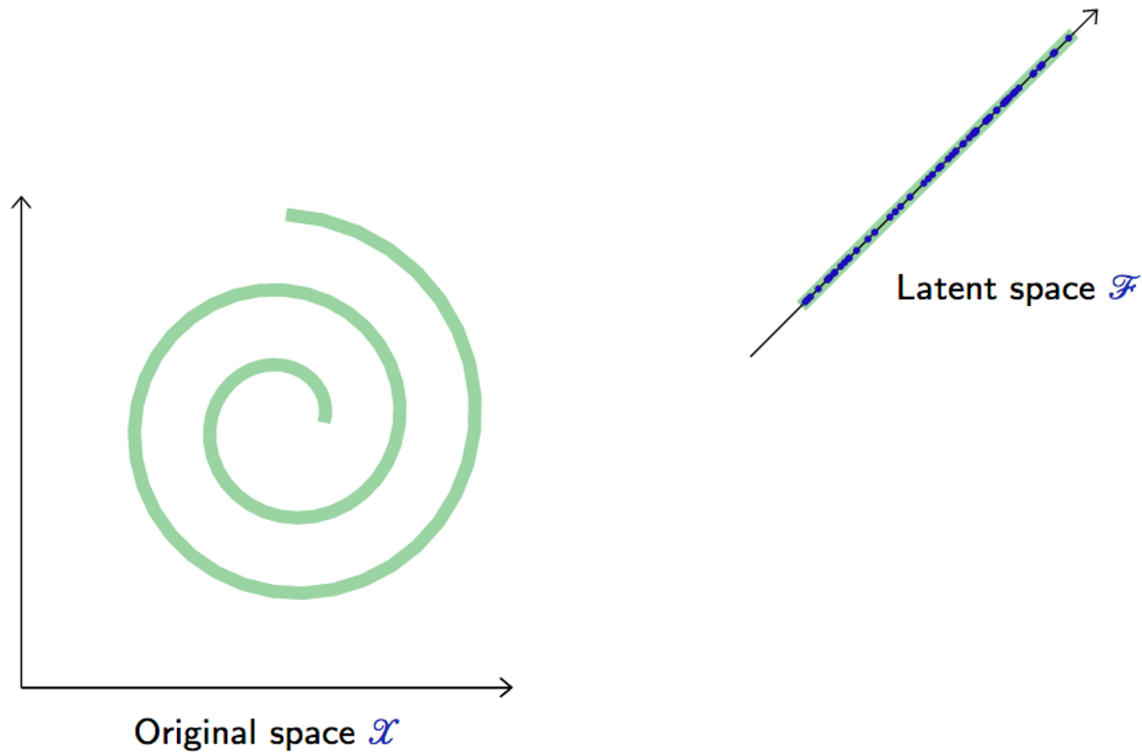
- Lecture 1
  - Introduction to Machine Learning fundamentals
  - Linear Models

- Lecture 2
  - Neural Networks
  - Deep Neural Networks
  - Convolutional, Recurrent, and Graph Neural Networks

- Lecture 3
  - Unsupervised Learning
  - Autoencoders
  - Generative Adversarial Networks and Normalizing Flows

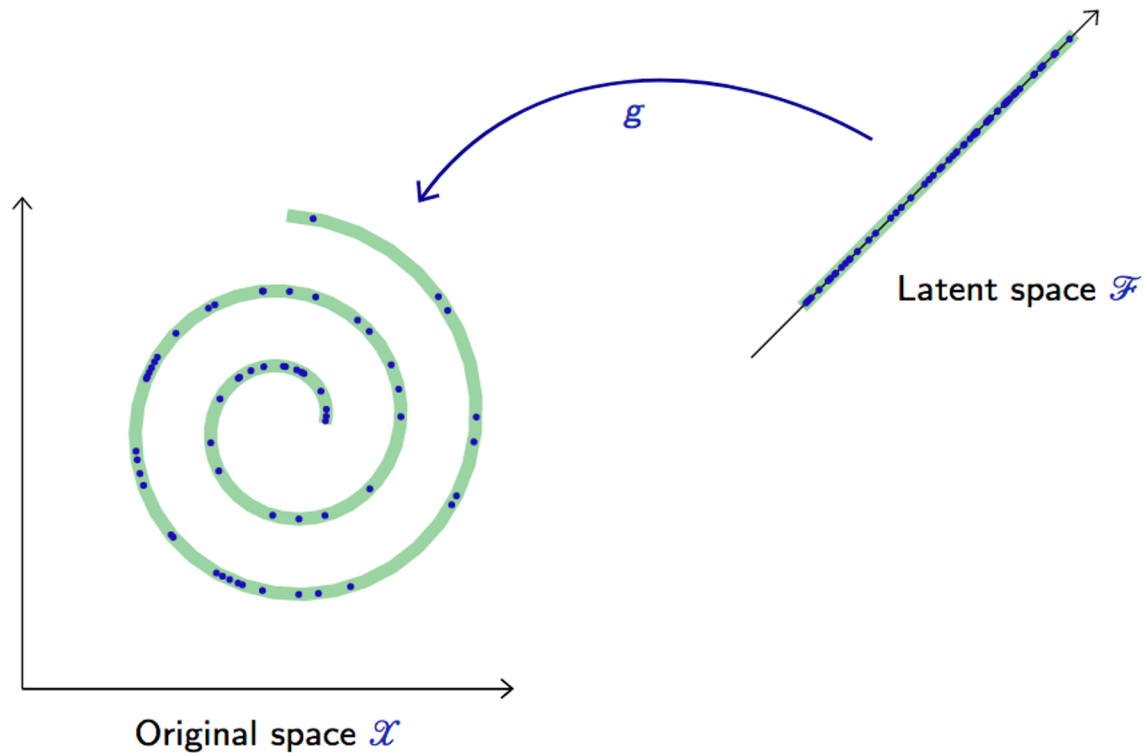# Beyond Regression and Classification

# Beyond Regression and Classification

- Not all tasks are predicting a label from features, as in classification and regression

- May want / need to explicitly model a high-dim. signal
  - Data synthesis / simulation
  - Density estimation
  - Anomaly detection
  - Denoising, super resolution
  - Data compression
  - …

- Often don't have labels → Unsupervised Learning

- Often framed as **modeling the lower dimensional "meaningful degrees of freedom"** that describe the data

Original space $\mathcal{X}$

Fleuret, Deep Learning Course

$f$

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

Original space $\mathcal{X}$

Latent space $\mathcal{F}$

Fleuret, Deep Learning Course

Original space $\mathcal{X}$

Latent space $\mathcal{F}$

$g$

# Autoencoders

# Meaningful Representations

- How can we find the "meaningful degrees of freedom" in the data?

- Dimensionality Reduction / Compression
  - Can we compress the data to a *latent space* with smaller number of dimensions, and still recover the original data from this latent space representation?

  - Latent space must encode and retain the important information about the data

  - Can we learn this compression and latent space

# Autoencoders

- Autoencoders map a space to itself through a compression, $x \to z \to \hat{x}$, and should be close to the identity on the data

  - Data: $x \in \mathcal{X}$         Latent space: $z \in \mathcal{F}$

  - **Encoder**: Map from $\mathcal{X}$ to a lower dimensional latent space $\mathcal{F}$
    - Parameterize as neural network $f_\theta(x)$ with parameters $\theta$

  - **Decoder**: Map from latent space $\mathcal{F}$ back to data space $\mathcal{X}$
    - Parameterize as neural network $g_\psi(z)$ with parameters $\psi$

# Autoencoders

- Autoencoders map a space to itself through a compression, $x \rightarrow z \rightarrow \hat{x}$, and should be close to the identity on the data

  - Data: $x \in \mathcal{X}$        Latent space: $z \in \mathcal{F}$

  - **Encoder**: Map from $\mathcal{X}$ to a lower dimensional latent space $\mathcal{F}$
    - Parameterize as neural network $f_\theta(x)$ with parameters $\theta$

  - **Decoder**: Map from latent space $\mathcal{F}$ back to data space $\mathcal{X}$
    - Parameterize as neural network $g_\psi(z)$ with parameters $\psi$

- What is the latent space? What are $f(x)$ and $g(z)$?
  - Choose a latent space dimension D
  - Learn mappings $f(x)$ to representation of size D, and back with $g(z)$

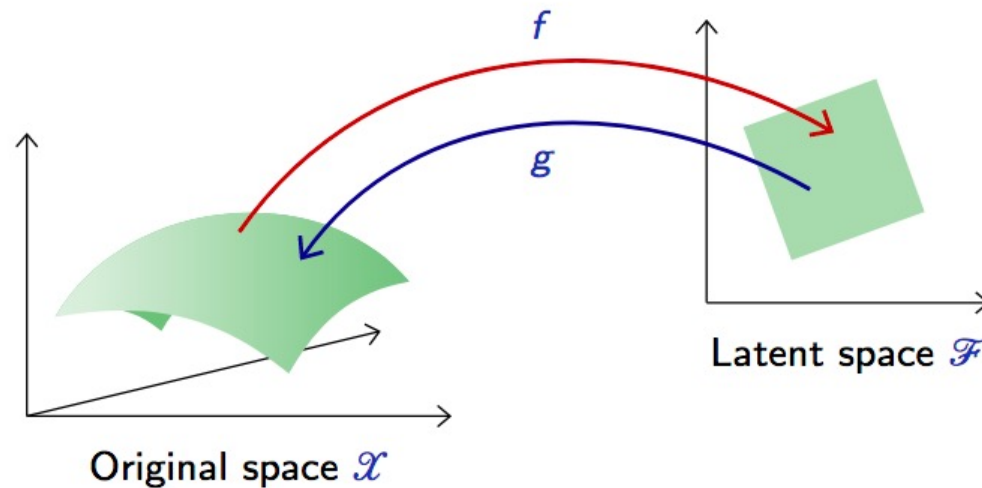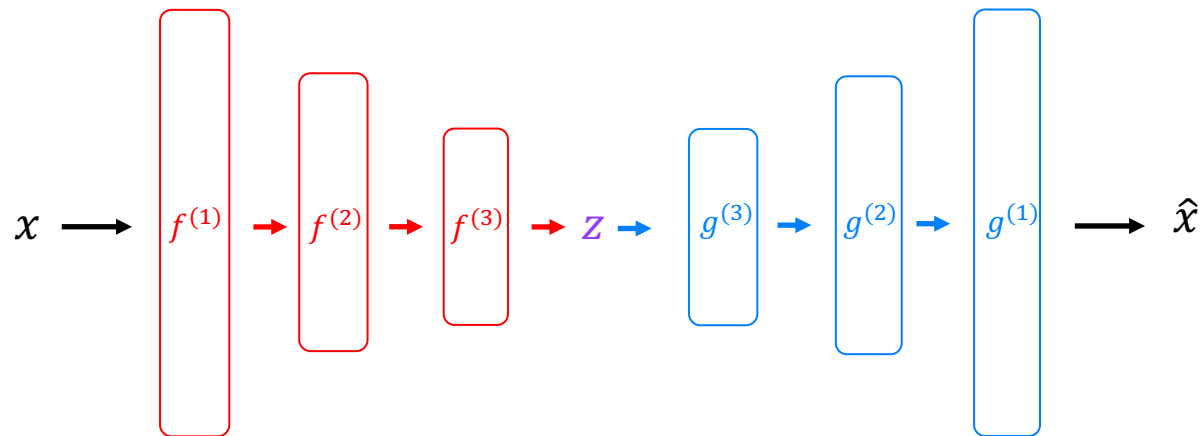- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data

$$L(\theta, \psi) = \frac{1}{N} \sum_n \left\| x_n - g_\psi\big(f_\theta(x_n)\big) \right\|^2$$

- Minimize this loss over parameters of encoder ($\theta$) and decoder ($\psi$).

- Loss: mean *reconstruction loss* (MSE) between data and encoded-decoded data

$$L(\textcolor{red}{\theta}, \textcolor{blue}{\psi}) = \frac{1}{N} \sum_n \left\| x_n - \textcolor{blue}{g_\psi}\big(\textcolor{red}{f_\theta}(x_n)\big) \right\|^2$$

- Minimize this loss over parameters of encoder ($\textcolor{red}{\theta}$) and decoder ($\textcolor{blue}{\psi}$).

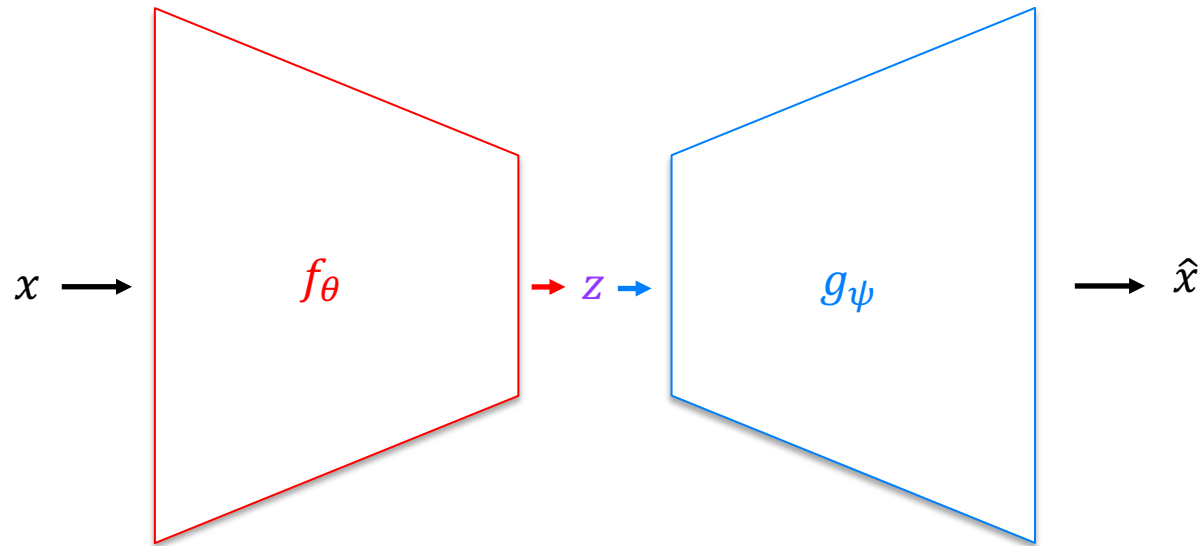- NOTE: if $f_\theta(x)$ and $g_\psi(z)$ are linear, optimal solution given by Principle Components Analysis

*f*

*g*

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

- If the latent space is of lower dimension, the autoencoder has to capture a "good" parametrization, and in particular dependencies between components

# Deep Autoencoder

$x \longrightarrow$ $f^{(1)}$ → $f^{(2)}$ → $f^{(3)}$ → $z$ → $g^{(3)}$ → $g^{(2)}$ → $g^{(1)}$ $\longrightarrow$ $\hat{x}$

- When $f_\theta$ and $g_\psi$ are multiple neural network layers, can learn complex mappings between $\mathcal{X}$ and $\mathcal{F}$
  - $f_\theta$ and $g_\psi$ can be Fully Connected, CNNs, RNNs, etc.
  - Choice of network structure will depend on data

- When $f_\theta$ and $g_\psi$ are multiple neural network layers, can learn complex mappings between $\mathcal{X}$ and $\mathcal{F}$
  - $f_\theta$ and $g_\psi$ can be Fully Connected, CNNs, RNNs, etc.
  - Choice of network structure will depend on data

# Deep Convolutional Autoencoder

$X$ (original samples)



$g \circ f(X)$ (CNN, $d = 16$)



$f_\theta$ and $g_\psi$ are each
5 convolutional layers

$g \circ f(X)$ (PCA, $d = 16$)



Fleuret, Deep Learning Course

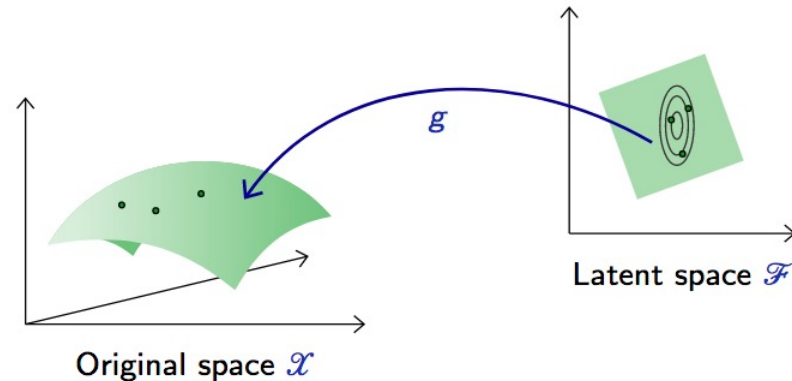$$\alpha \in [0,1], \quad \xi(x, x', \alpha) = g((1-\alpha)f(x) + \alpha f(x')).$$



Autoencoder interpolation ($d = 8$)

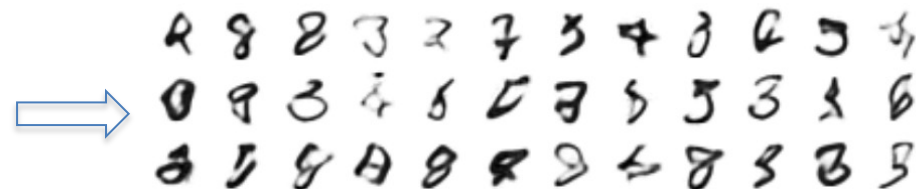- Can we sample in latent space and decode to generate data?

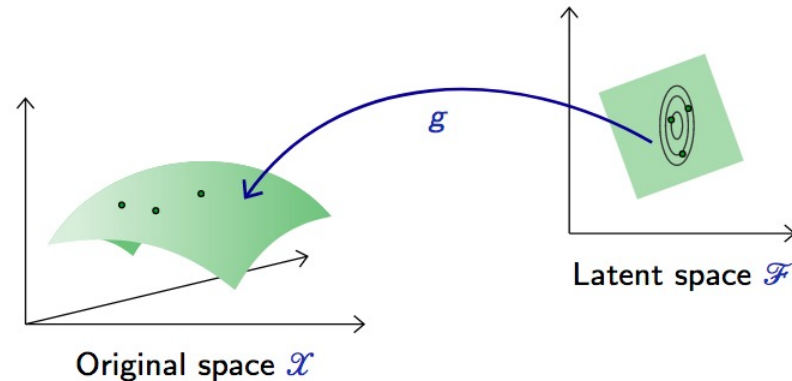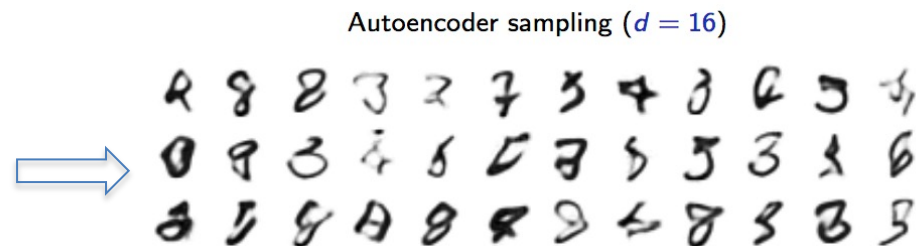- Can we sample in latent space and decode to generate data?



- What distribution to sample from in latent space?
  - Try Gaussian with mean and variance from data

- Can we sample in latent space and decode to generate data?



- What distribution to sample from in latent space?
  - Try Gaussian with mean and variance from data



- Don't know the right latent space density

# Generative Models

# Generative Models

- Generative models aim to:
  - Learn a distribution $p(x)$ that explains the data
  - Draw samples of plausible data points

- Explicit Models
  - Can evaluate the density $p(x)$ of a data point x

- Implicit Models
  - Can only sample from $p(x)$, but not evaluate density

- Learn a mapping from corrupted data space $\widetilde{\mathcal{X}}$ back to original data space

  - Mapping $\phi_w(\widetilde{x}) = x$
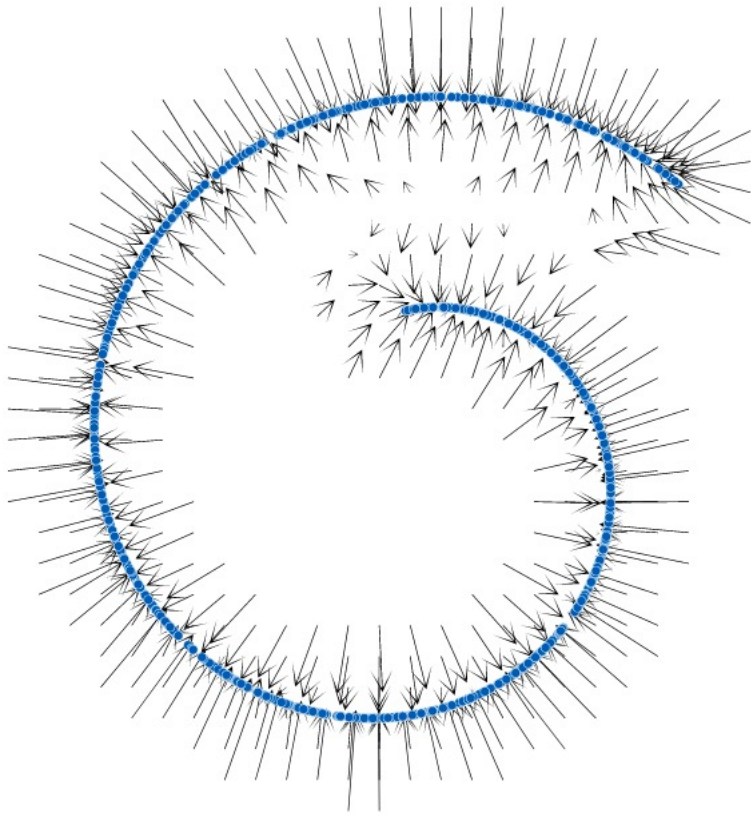  - $\phi_w$ will be a neural network with parameters $w$

- Loss:

$$L = \frac{1}{N} \sum_n \|x_n - \phi_w(x_n + \epsilon_n)\|$$

Perturbation, e.g. Gaussian noise
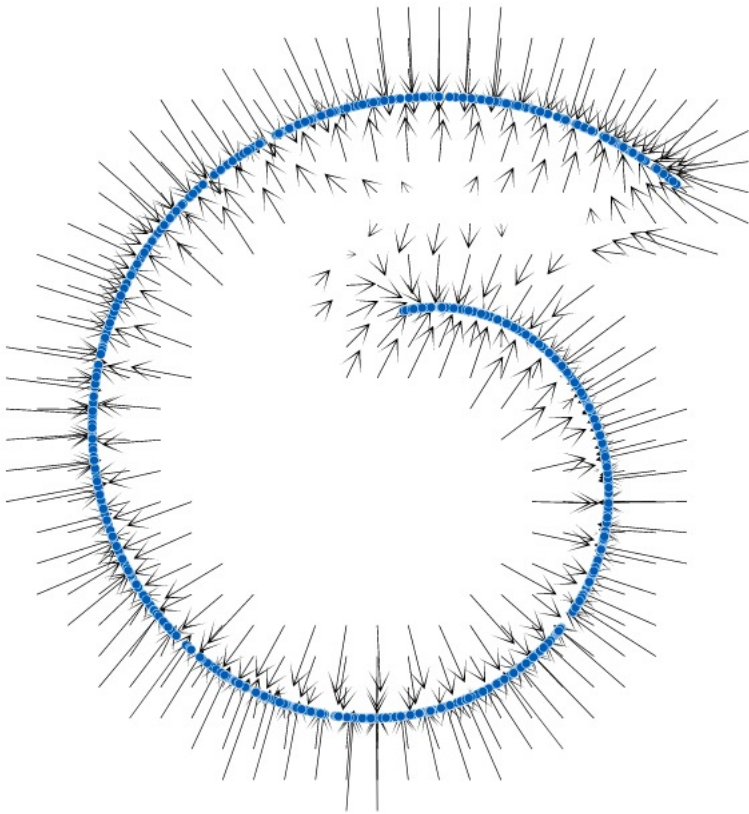
# Denoising Autoencoders Examples

Original

Corrupted ($\sigma = 4$)

Reconstructed

- Autoencoder learns the average behavior

- What if we care about these variations?
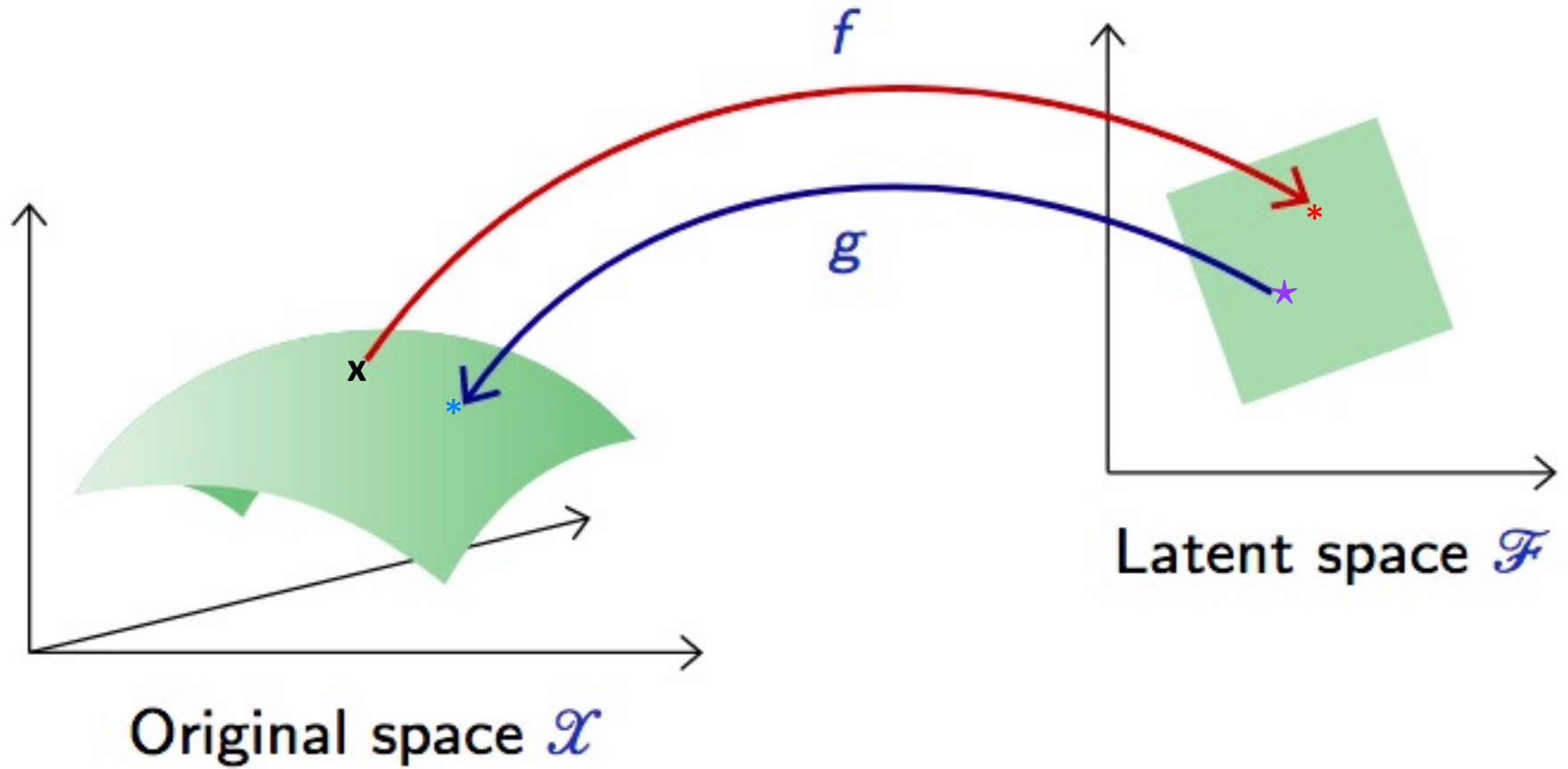
- Can we add a notion of variation in the autoencoder?

- Consider probabilistic relationship between data and latent variables
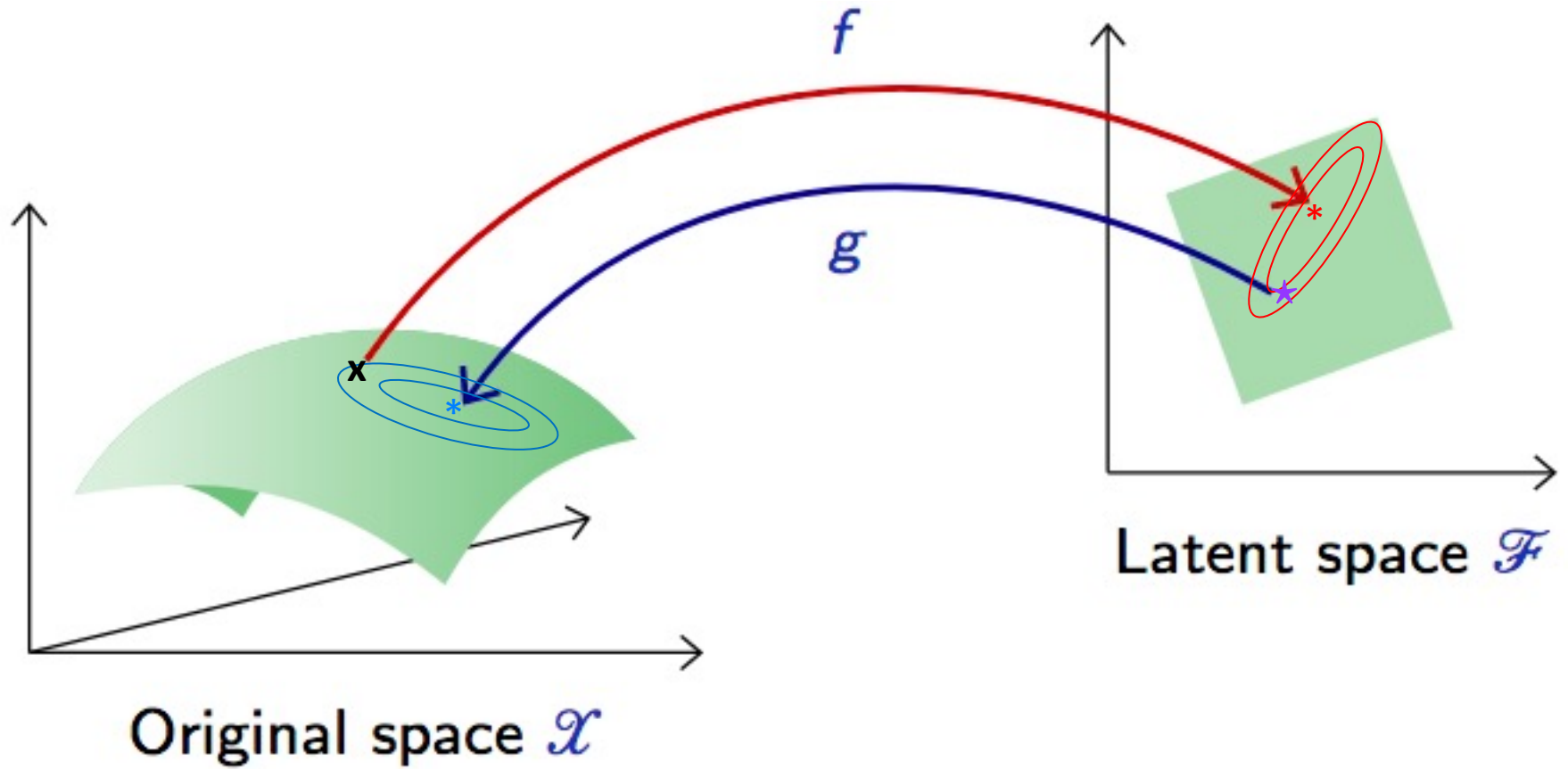
$$x, z \sim p(x, z) = p(x|z)p(z)$$

- Autoencoding

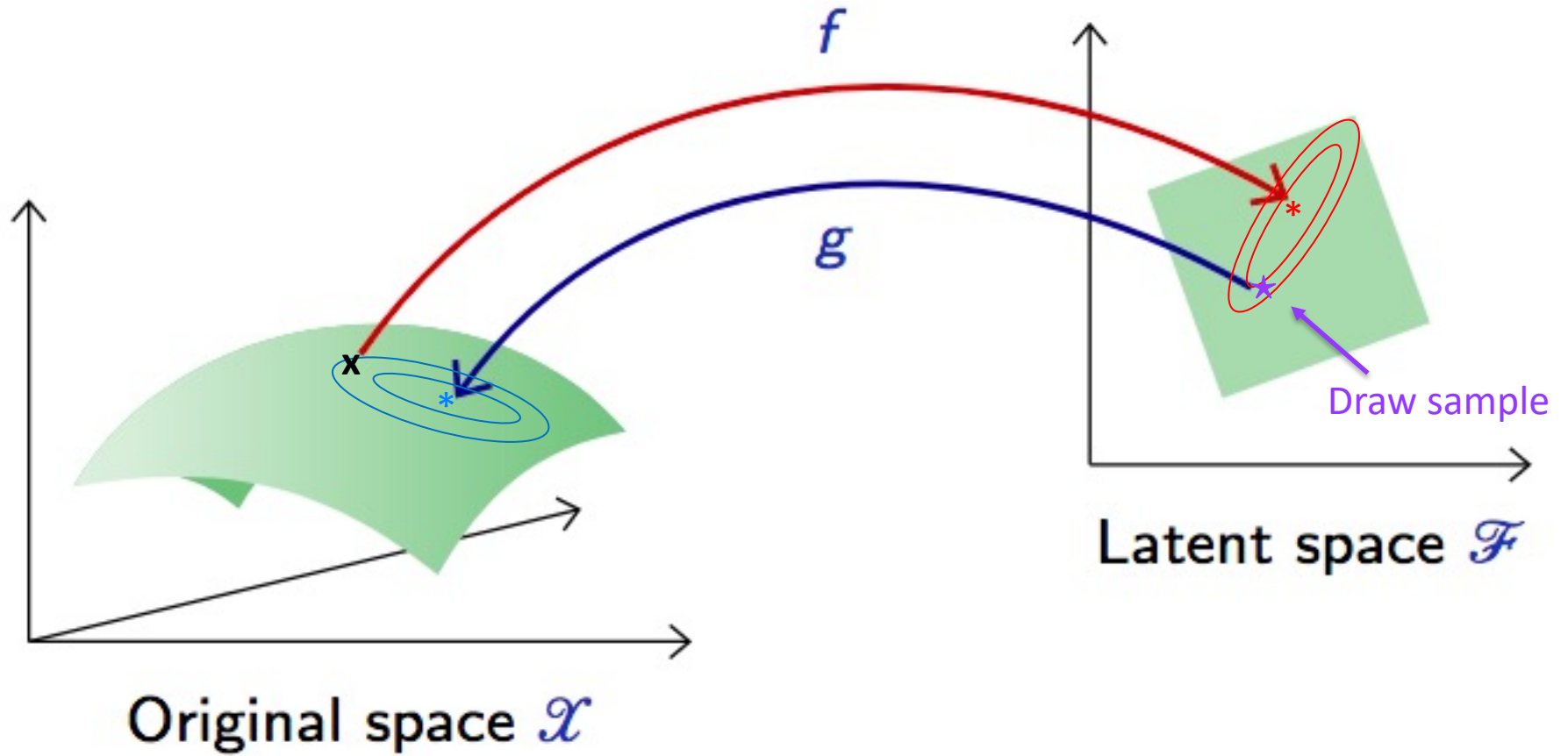$$x \rightarrow q(z|x) \underset{sample}{\Longrightarrow} z \rightarrow p(x|z)$$

  – Choose simple prior distribution

  – **Encoder:** Learn what latents can produced data: $q(z|x)$
  – **Decoder:** Learn what data is produced by latent: $p(x|z)$

$f$

$g$

Latent space $\mathscr{F}$

x

Original space $\mathscr{X}$

$f$

$g$

x

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

Original space $\mathcal{X}$

Latent space $\mathcal{F}$

Draw sample
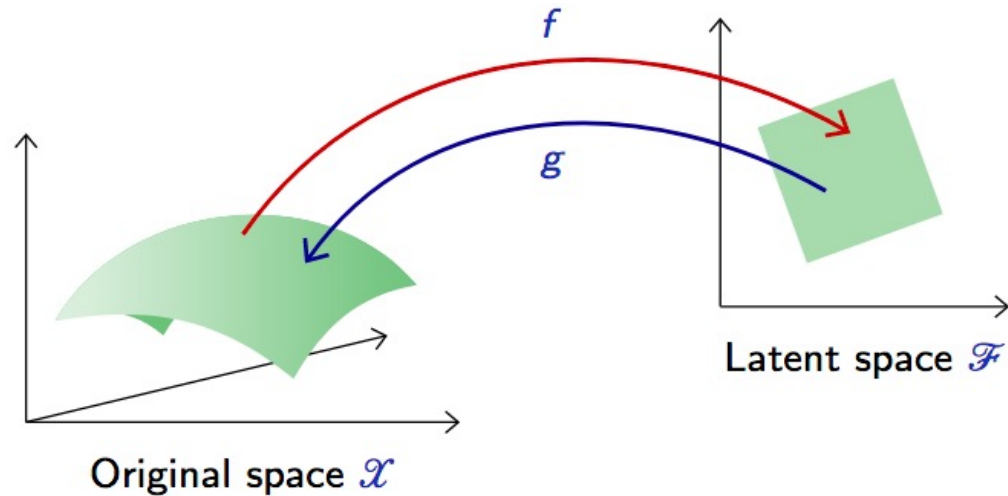
- Typical encoder maps input $x$ to "average" point in latent space

$$f(x) = \mu(x)$$

- A VAE Encoder has two outputs: mean & variance function

$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi(x)\}$$

$\psi$ are parameters of the NN

# Encoding

- A VAE Encoder has two outputs: mean & variance function

$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi(x)\}$$   $\psi$ are parameters of the NN

- What is the probability of a point in latent space?

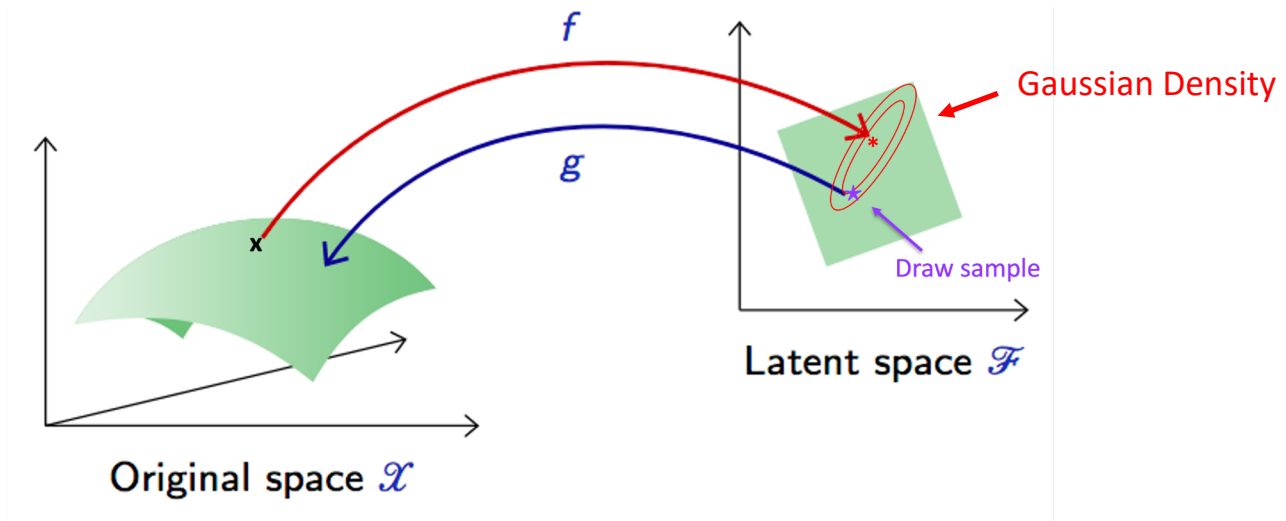$$p_\psi(z|x) = N(z \mid \mu_\psi(x), \sigma_\psi(x))$$

# Encoding

- A VAE Encoder has two outputs: mean & variance function

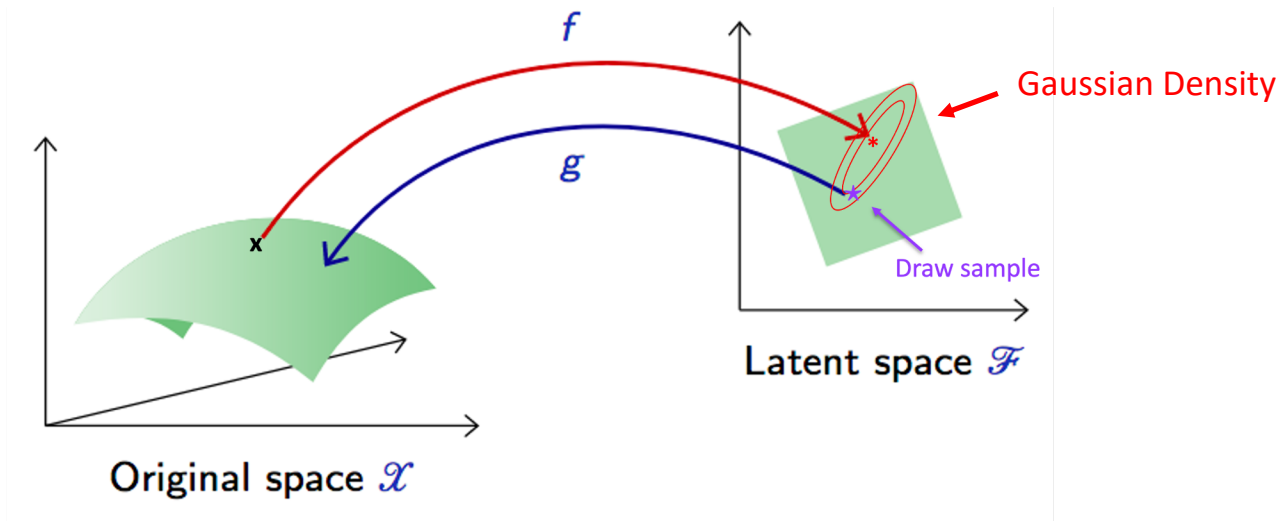$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi(x)\}$$

$\psi$ are parameters of the NN

- What is the probability of a point in latent space?

$$p_\psi(z|x) = N(z \mid \mu_\psi(x), \sigma_\psi(x))$$

- How do we draw a sample in latent space?

$$z = \sigma_\psi(x) * \epsilon + \mu_\psi(x) \qquad \epsilon \sim N(0, I)$$

Re-parameterization trick



Gaussian Density

Draw sample

Latent space $\mathcal{F}$

Original space $\mathcal{X}$
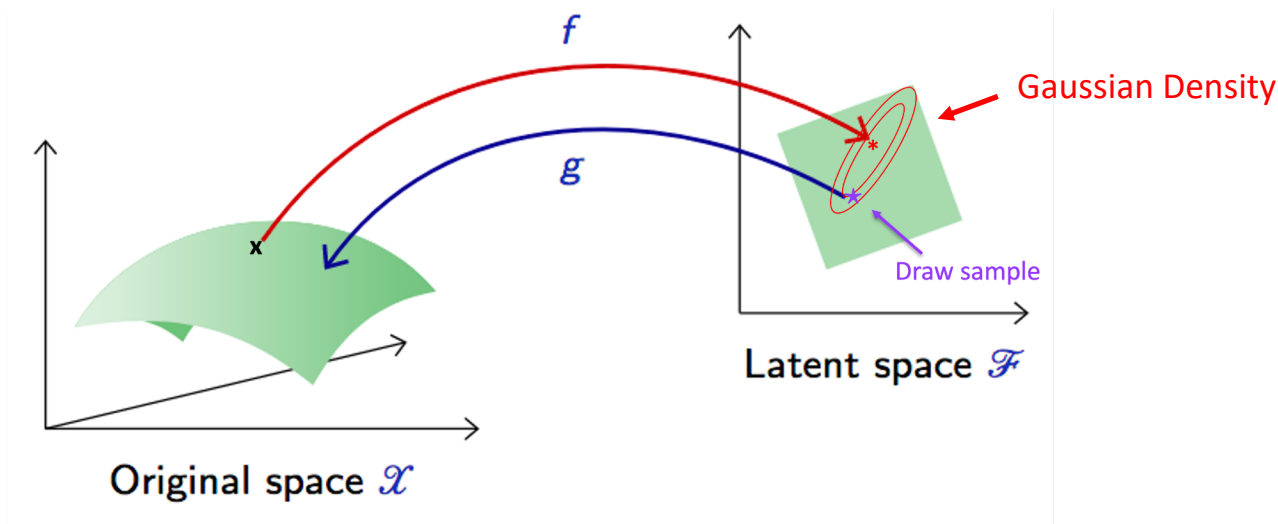
# Encoding

- A VAE Encoder has two outputs: mean & variance function

$$f_\psi(x) = \{\mu_\psi(x), \sigma_\psi(x)\}$$

$\psi$ are parameters of the NN
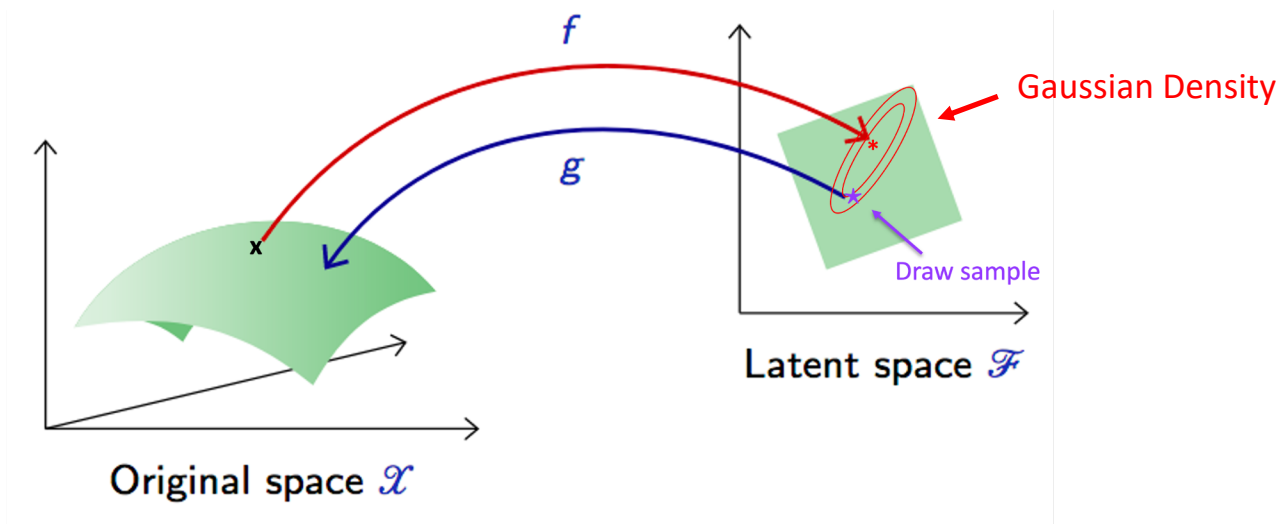
- What is the probability of a point in latent space?

$$p_\psi(z|x) = N(z \mid \mu_\psi(x), \sigma_\psi(x))$$

- How do we draw a sample in latent space?

$$z = \sigma_\psi(x) * \epsilon + \mu_\psi(x) \qquad \epsilon \sim N(0, I)$$

Re-parameterization trick

NOTE:
Could have chosen
different density and
use NN to predict
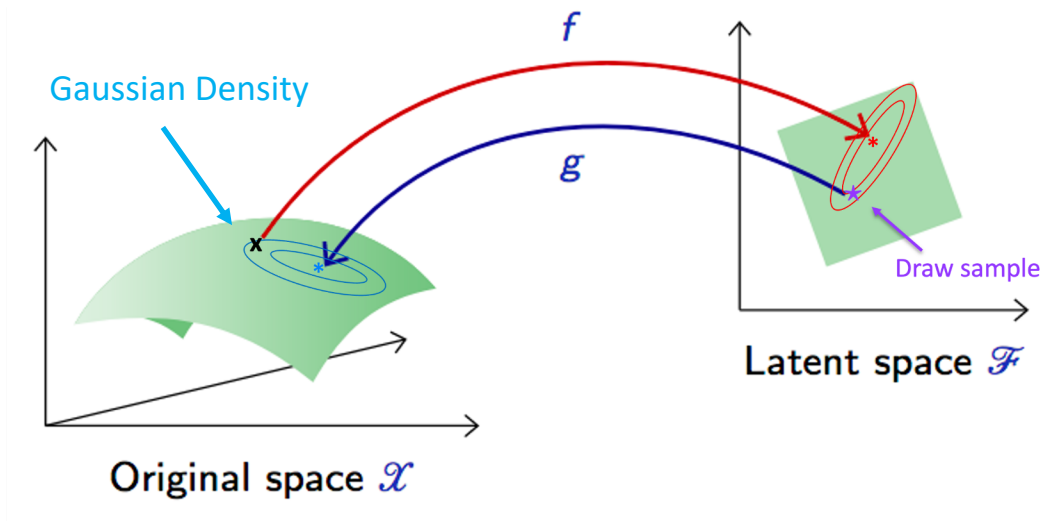params…

As long as we can
sample using
re-parameterization



Gaussian Density

Draw sample

Latent space $\mathcal{F}$

Original space $\mathcal{X}$

# Decoding
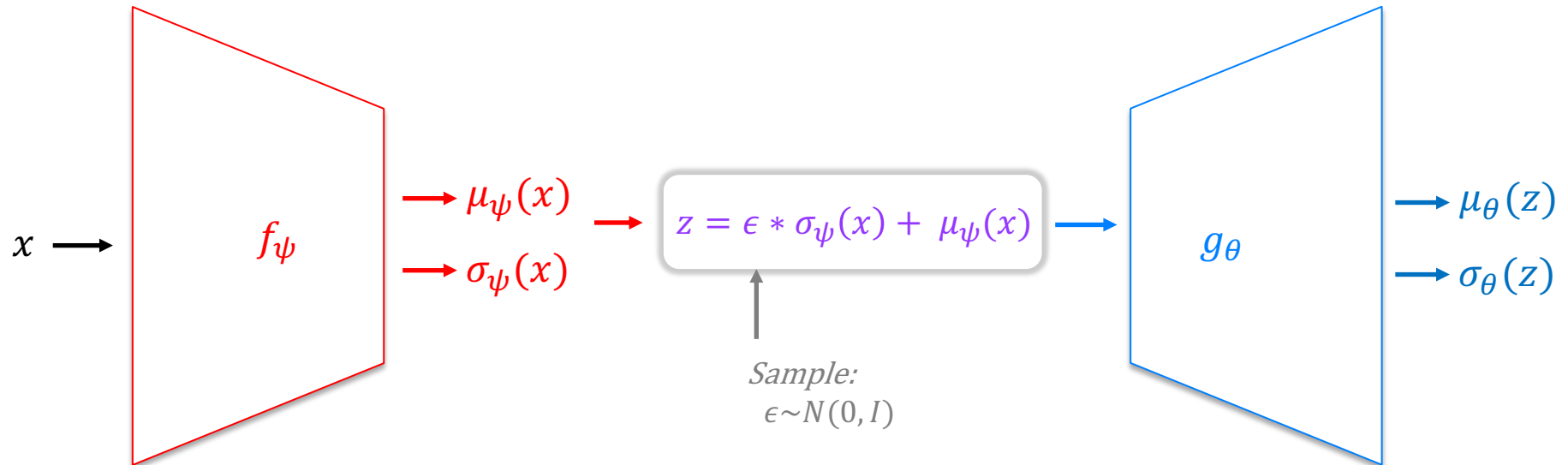
- Same approach, VAE decoder has two outputs

$$g_\theta(z) = \{\mu_\theta(z), \sigma_\theta(z)\}$$

$\theta$ are parameters of the NN

- Likelihood of an observation $x$

$$p_\theta(x|z) = N(x \mid \mu_\theta(z), \sigma_\theta(z))$$



Gaussian Density

*f*

*g*

Draw sample

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

$x$

$f_\psi$

$\mu_\psi(x)$

$\sigma_\psi(x)$

$z = \epsilon * \sigma_\psi(x) + \mu_\psi(x)$

*Sample:*
$\epsilon \sim N(0, I)$

$g_\theta$

$\mu_\theta(z)$

$\sigma_\theta(z)$

*f*

*g*

Draw sample

Latent space $\mathscr{F}$

Original space $\mathscr{X}$

***Reconstruction Loss***: Maximize expected likelihood of decoding $x$ from encodings of $x$

$$L_{reco} = \mathbb{E}_{z \sim q(z|x)}[\log p(x|z)] \approx \frac{1}{N} \sum_{z_i \sim q(z|x)} \log p(x|z_i)$$

# Variational Autoencoder Training Loss

- $L_{reco} = \frac{1}{N} \sum_{z \sim q_\psi(z|x)} \log p_\theta(x|z_i)$

- Note that

$$\log p(x|z) = -\log \sigma_\theta(z) - \frac{\left(x - \mu_\theta(z)\right)^2}{\sigma_\theta(z)^2} + const$$

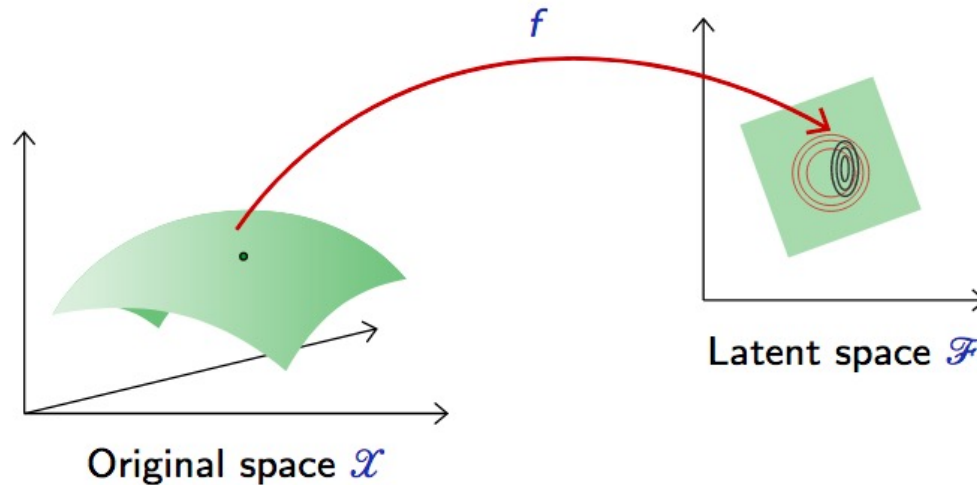This looks almost exactly like the Autoencoder Loss

Which was a Mean Squared Error $\left(x - f\left(g(x)\right)\right)^2$

Here we have $z \equiv z_\psi(x)$

- $L_{reco} = \frac{1}{N} \sum_{z \sim q_\psi(z|x)} \log p_\theta(x|z_i)$

- What about encoder? How do we make sure it doesn't collapse around each point (i.e. only predict mean)

- $L_{reco} = \frac{1}{N} \sum_{z \sim q_\psi(z|x)} \log p_\theta(x|z_i)$

- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**

- $L_{reco} = \frac{1}{N} \sum_{z \sim \color{red}{q_\psi(z|x)}} \log \color{blue}{p_\theta(x|z_i)}$

- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**

- Constrain difference between distributions with **Kullback–Leibler divergence**

$$D_{KL}[q(z|x)|p(z)] = \mathbb{E}_{q(z|x)}\left[\log \frac{q(z|x)}{p(z)}\right] = \int q(z|x) \log \frac{q(z|x)}{p(z)} \, dz$$
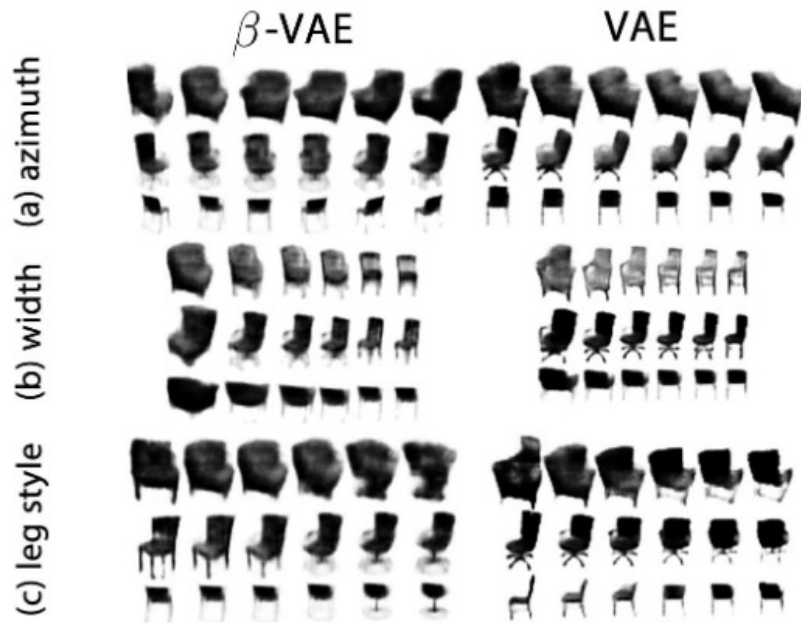
  – $D_{KL}[q|p] \geq 0$  and is only 0 when $q = p$
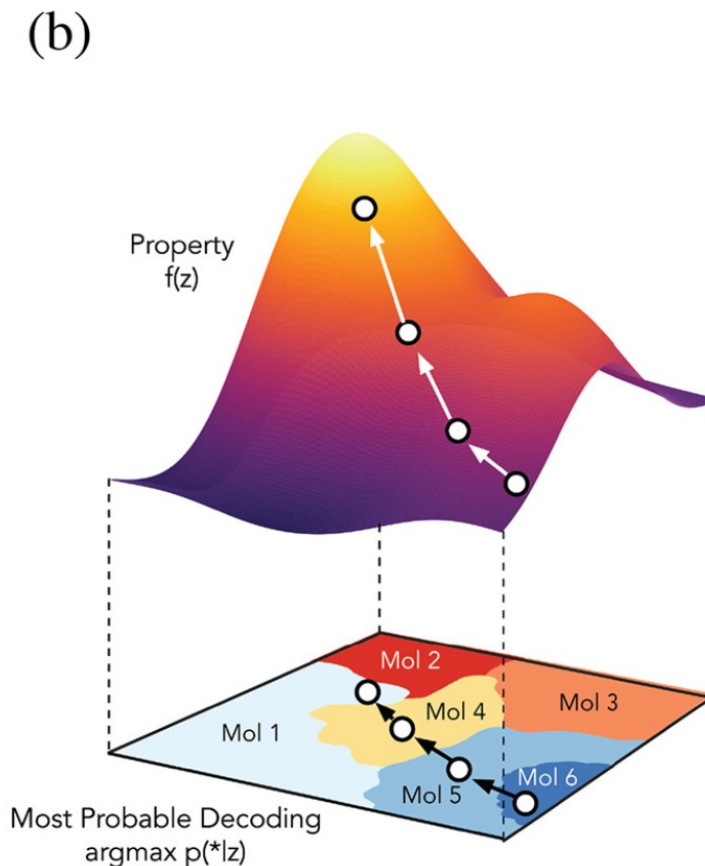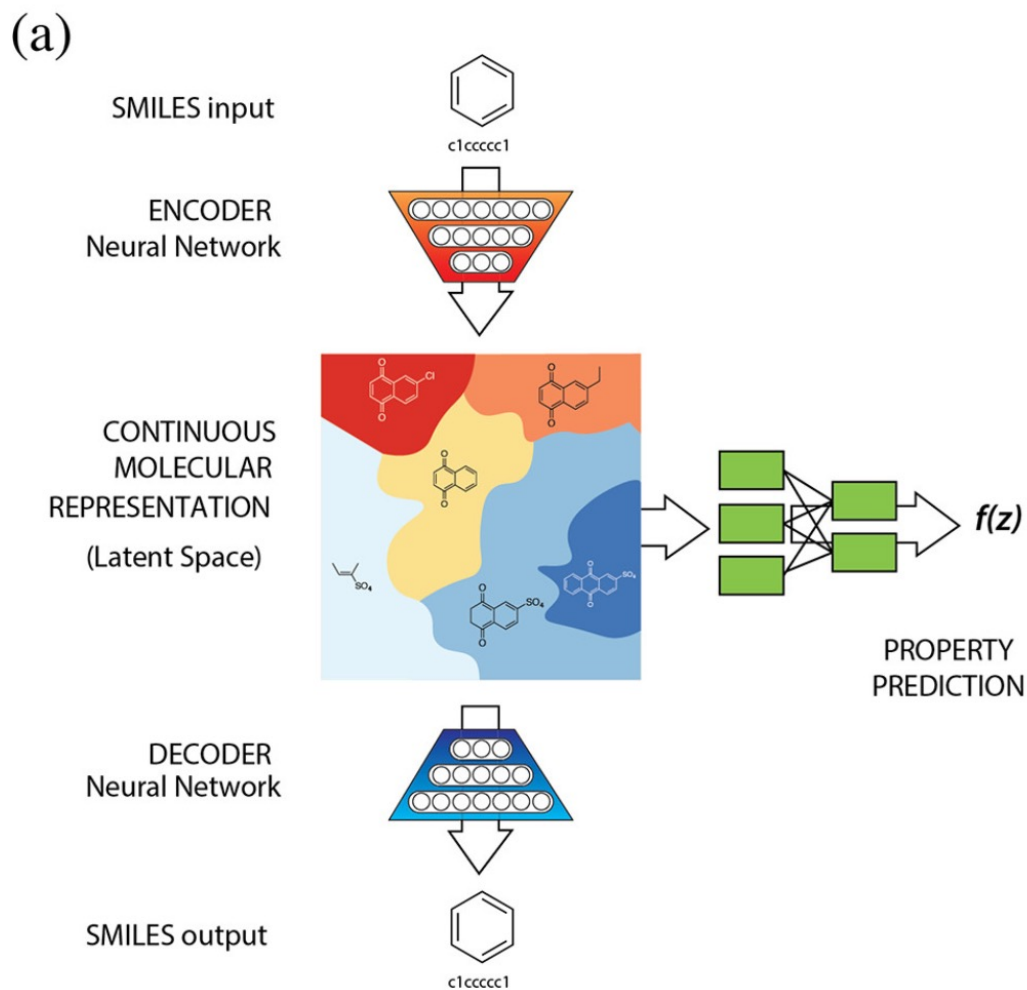
# Variational Autoencoder Training Loss

- $L_{reco} = \frac{1}{N}\sum_{z \sim q_\psi(z|x)} \log p_\theta(x|z_i)$

- Use prior $p(z)$ for the latent space distribution, **need to ensure the encoder is consistent with prior**

- VAE full objective

$$\max_{\theta,\psi} L(\theta,\psi) = \max_{\theta,\psi} \left[ \mathbb{E}_{q_\psi(z|x)}[\log p_\theta(x|z)] - D_{KL}[q_\psi(z|x)|p(z)] \right]$$

Reconstruction Loss          Regularization of Encoder

Higgins et al., 2017

Design of new molecules with desired chemical properties.
(Gomez-Bombarelli et al, 2016)

# Another Way To Do Generative Modeling…

- Formulate as a two player game

- One player tries to output data that looks as real as possible

- Another player tries to compare real and fake data

- In this case we need:
    1. A *generator* that can produce samples
    2. A measure of  *not too far from the real data*

- **Generator network $g_\theta(z)$** with parameters $\theta$
  - Map sample from known $p(z)$ to sample in data space

$$x = g_\theta(z) \quad z \sim p(z)$$

  - We don't know what the generated distribution $p_\theta(x)$ is, but we can sample from it → *Implicit Model*

# Generative Adversarial Network (GAN)

- **Generator network $g_\theta(z)$** with parameters $\theta$
  - Map sample from known $p(z)$ to sample in data space
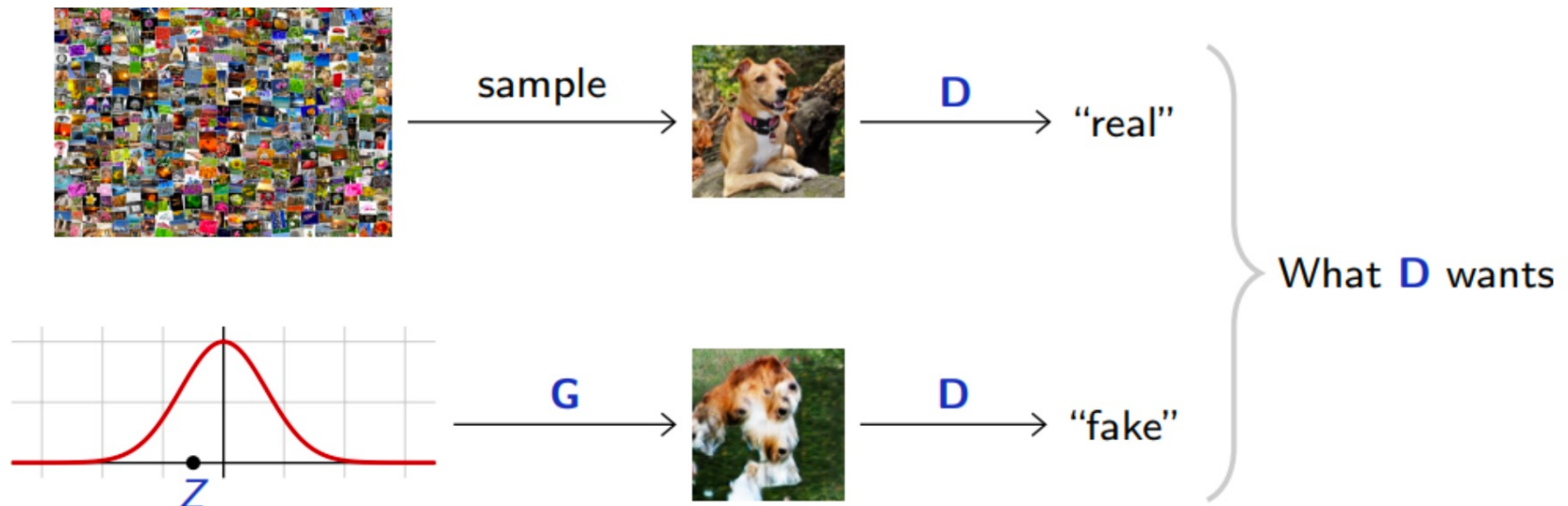
$$x = g_\theta(z) \quad z \sim p(z)$$

  - We don't know what the generated distribution $p_\theta(x)$ is, but we can sample from it → *Implicit Model*

- **Discriminator Network $d_\phi(x)$** with parameters $\phi$
  - Classifier trained to distinguish between real and fake data

  - Classifier is learning to predict $p(y = real \mid x)$

  - This classifier is our measure of *not too far from the real data*

# GAN Setup

- Generator's goal is to produce *fake* data that tricks the discriminator to think it is *real* data

- Discriminator wants to miss-classify data as real or fake as little as possible

- The setup is *adversarial* because the two networks have opposing objectives

# GAN Objective

- Data
  - Real data samples: $\{x_i, y_i = 1\}$

  - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$  with: $z_i \sim p(z)$

# GAN Objective

- Data
  - Real data samples: $\{x_i, y_i = 1\}$

  - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$ with: $z_i \sim p(z)$

- For a fixed generator, can train discriminator by minimizing the cross entropy

$$L(\phi) = -\frac{1}{2N} \sum_{i=1}^{N} \left[ y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i)) \right]$$

# GAN Objective

- Data
  - Real data samples: $\{x_i, y_i = 1\}$

  - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$    with: $z_i \sim p(z)$

- For a fixed generator, can train discriminator by minimizing the cross entropy

$$L(\phi) = -\frac{1}{2N} \sum_{i=1}^{N} \left[ y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i) ) \right]$$

$$= -\frac{1}{2N} \sum_{i=1}^{N} \left[ \log d_\phi(x_i) + \log(1 - d_\phi(g_\theta(z_i)) ) \right]$$

# GAN Objective

- Data
  - Real data samples: $\{x_i, y_i = 1\}$

  - Fake data samples: $\{\tilde{x}_i = g_\theta(z_i), \tilde{y}_i = 0\}$  with: $z_i \sim p(z)$

- For a fixed generator, can train discriminator by minimizing the cross entropy

$$L(\phi) = -\frac{1}{2N} \sum_{i=1}^{N} \left[ y_i \log d_\phi(x_i) + (1 - \tilde{y}_i) \log(1 - d_\phi(\tilde{x}_i)) \right]$$

$$= -\frac{1}{2N} \sum_{i=1}^{N} \left[ \log d_\phi(x_i) + \log(1 - d_\phi(g_\theta(z_i))) \right]$$

$$= -\mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ \log d_\phi(x) \right] - \mathbb{E}_{z \sim p(z)} \left[ \log(1 - d_\phi(g_\theta(z))) \right]$$

# GAN Objective

- However, generator isn't fixed… have to train it!

- However, generator isn't fixed… have to train it!

- Consider objective as a *value function* of $\phi$ and $\theta$

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \Big[ \log d_\phi(x) \Big] + \mathbb{E}_{z \sim p(z)} \Big[ \log(1 - d_\phi(g_\theta(z)) ) \Big]$$

- However, generator isn't fixed… have to train it!

- Consider objective as a *value function* of $\phi$ and $\theta$

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \left[ \log d_\phi(x) \right] + \mathbb{E}_{z \sim p(z)} \left[ \log(1 - d_\phi(g_\theta(z))\, ) \right]$$

  – For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes

  – For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low

# GAN Objective

- However, generator isn't fixed… have to train it!

- Consider objective as a *value function* of $\phi$ and $\theta$

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)} \Big[ \log d_\phi(x) \Big] + \mathbb{E}_{z \sim p(z)} \Big[ \log(1 - d_\phi(g_\theta(z)) ) \Big]$$

  – For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes

  – For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low

- So our optimization goal becomes:
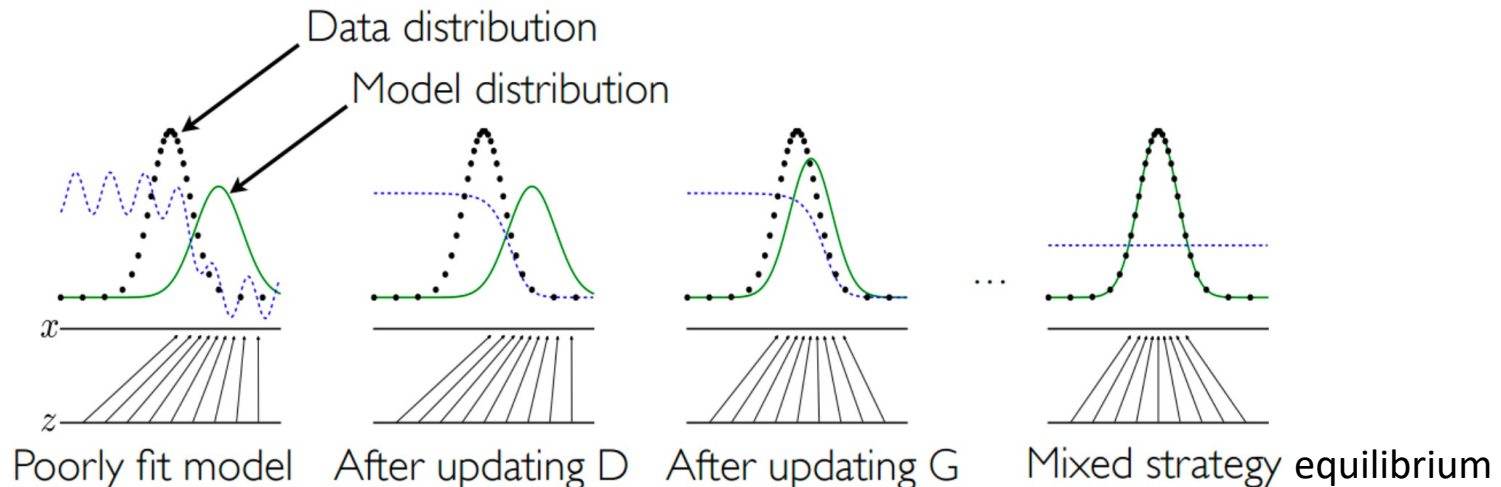
$$\theta^* = \arg \min_\theta \max_\phi V(\phi, \theta)$$

# GAN Objective

- However, generator isn't fixed… have to train it!

- Consider objective as a *value function* of $\phi$ and $\theta$

$$V(\phi, \theta) = \mathbb{E}_{x \sim p_{\text{data}}(x)}\Big[\log d_\phi(x)\Big] + \mathbb{E}_{z \sim p(z)}\Big[\log(1 - d_\phi(g_\theta(z)) )\Big]$$

  - For fixed generator, $V(\phi, \theta)$ is high when discriminator is good, i.e. when generator is not producing good fakes

  - For a perfect discriminator, a good generator will confuse discriminator and $V(\phi, \theta)$ will be low

- So our optimization goal becomes:

$$\theta^* = \arg \min_\theta \max_\phi V(\phi, \theta)$$

NOTE: can prove that minimax solution corresponds to generator that perfectly reproduces data distribution
$q_{\theta^*}(x) = p_{data}(x)$

# GAN Training

- Alternating Gradient descent to solve the min-max problem:

$$\theta \leftarrow \theta - \gamma \nabla_\theta V(\phi, \theta) = \theta - \gamma \frac{\partial V}{\partial d} \frac{\partial (d_\phi)}{\partial g} \frac{\partial g_\theta}{\partial \theta}$$

$$\phi \leftarrow \phi - \gamma \nabla_\phi V(\phi, \theta) = \phi - \gamma \frac{\partial V}{\partial d} \frac{d(d_\phi)}{d\phi}$$

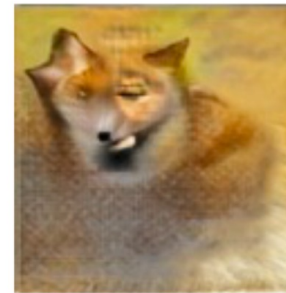- For each $\theta$ step, take $k$ steps in $\phi$ to keep discriminator near optimal



Data distribution
Model distribution

$x$

$z$

Poorly fit model    After updating D    After updating G    Mixed strategy equilibrium

[GAN Lab Demo](#)

# Examples

Goodfellow et. al., 2014

Not so good

Goodfellow 2016

Radford et al, 2015

# Challenges

- **Oscillations without convergence**: unlike standard loss minimization, alternating stochastic gradient descent has no guarantee of convergence.

- **Vanishing gradients**: if classifier is too good, value function saturates → no gradient to update generator

- **Mode collapse**: generator models only a small sub-population, concentrating on a few data distribution modes.

- **Difficult to assess performance**, when are generated data good enough?



Step 0    Step 5k    Step 10k    Step 15k    Step 20k    Step 25k    Target

Slide credit: G. Louppe                                    Mode collapse (Metz et al, 2016)

- Standard GANS compare real and fake distributions with Jensen–Shannon Divergence, "vertically"

- Wasserstein-GAN (Arjovsky et al, 2017) compares "horizontally" with Wasserstein–1 distance (a.k.a. Earth Movers distance)

- Substantially improves *vanishing gradient* and *mode collapse* problems!

(Arjovsky et al, 2017)

Density of real
Density of fake
GAN Discriminator
WGAN Critic

Linear gradients in a WGAN

Vanishing gradients in regular GAN

*Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.*

(Arjovsky et al, 2017)

Progressive GAN



(Karras et al, 2017)

# Scaling Up

StyleGAN v2



(Karras et al, 2019)

BigGAN



(Brock et al, 2018)

- $p(z)$ doesn't have to be random noise

- CycleGAN uses *cycle-consistency loss* in addition to GAN loss
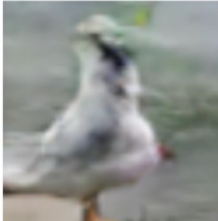  - Translating from A$\rightarrow$B$\rightarrow$A should be consistent with original A

Fig. 3: Example results by our StackGAN-v1, GAWWN [29], and GAN-INT-CLS [31] conditioned on text descriptions from CUB test set.

(Zhang et al, 2017)

# Summary

- Deep neural networks are an extremely powerful class of models

- We can express our inductive bias about a system in terms of model design, and can be adapted to a many types of data

- Even beyond classification and regression, deep neural networks allow for powerful model schemes such as Variational Autoencoder and Generative adversarial Networks

- Must first determine the question we want to ask, and formulate an appropriate loss function
  - Loss function encodes the quality of model prediction
  - Parameterize models with neural networks

- Will have many of the same theoretical and practical issues as in classification and regression
  - What is the right class and structure of the model (CNN, RNN, graph, etc.) for the data?
  - How do we stably optimize the loss w.r.t. parameters?

- Autoencoders learn the latent space, but we don't know what is the latent space distribution

- Autoencoder prescribes a deterministic relationship between data space and latent space

- One set of "meaningful degrees of freedom" can only describe one data space point

- Observed random variable $x$ depends on unobserved latent random variable $z$
  - Interpret $z$ as the causal factors for $x$

- Joint probability: $p(x, z) = p(x|z)p(z)$

- $p(x|z)$ is a stochastic generation process from $z \rightarrow x$

- Inference from posterior: $\quad p(z|x) = \dfrac{p(x|z)p(z)}{p(x)}$

  - Usually can't compute marginal $p(x) = \int p(x|z)p(z)dz$

- Consider probabilistic relationship between data and latent variables

$$x, z \sim p(x, z) = p(x|z)p(z)$$
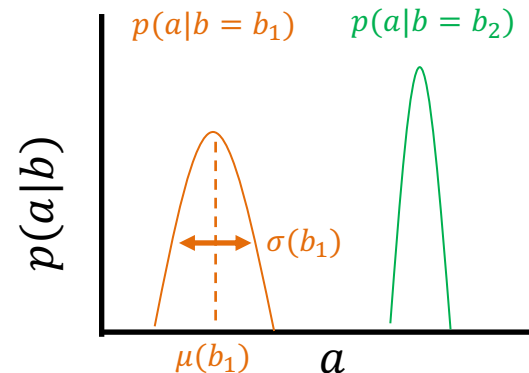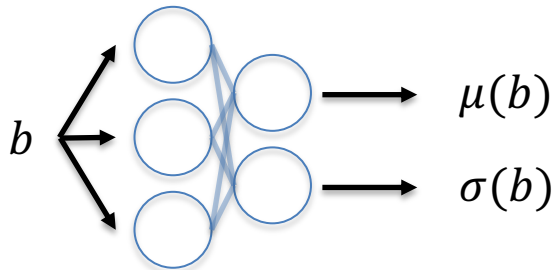
Decoding data x from latent z

Prior over latent space

- Classification / regression models make single predictions…

  How to model a conditional density $p(a|b)$ ?

- Assume a known form of density, e.g. normal

$$p(a|b) = \mathcal{N}\big(a; \mu(b), \sigma(b)\big)$$

  – Parameters of density depend on conditioned variable

- Use neural network to model density parameters

# The Decoder

- **Decoder**
  - Neural network with parameters $\theta$
  - Input $z$ → output estimate of Gaussian $\mu_\theta(z)$ , $\sigma_\theta(z)$

- **Likelihood of a data point x**

$$\log p(x|z) = -\log \sigma_\theta(z) - \frac{\left(x - \mu_\theta(z)\right)^2}{\sigma_\theta(z)^2} + const$$

- **Encoder**
  - Neural network with parameters $\psi$
  - Input $x$ → outputs estimate of Gaussian $\mu_\psi(x)$ , $\sigma_\psi(x)$

- For reconstruction loss:
  - Need a value of $z$ to evaluate decoder!
  - Need to gradient through $z$ to encoder parameters

$$\max_{\theta,\psi} L(\theta,\psi) = \max_{\theta,\psi} \sum_{z_i \sim q_\psi(z|x)} \log p_\theta(x|z_i) - \log\left[\frac{q_\psi(z_i|x)}{p(z_i)}\right]$$

- For $z \sim p_\theta(z)$, rewrite $z$ as a function of a random variable $\epsilon$ whose distributions $p(\epsilon)$ does not depend on $\theta$

  – Gaussian Example:

$$z \sim \mathcal{N}(\mu, \sigma) \quad \rightarrow \quad z = \sigma * \epsilon + \mu \quad where \;\; \epsilon \sim \mathcal{N}(0,1)$$

- VAE Loss

$$\max_{\theta, \psi} L(\theta, \psi) = \max_{\theta, \psi} \sum_{\epsilon \sim p(\epsilon)} \log p_\theta\big(x \big| z_i = \epsilon * \sigma_\psi(x) + \mu_\psi(x)\big) - \log \left[ \frac{q_\psi(z_i|x)}{p(z_i)} \right]$$

# Explicit Density Estimation with Normalizing Flows

# Explicit Density Estimation

- In VAE and GAN we can learn to sample from the distribution…

- Is there a way to learn the explicit density $p(x)$ ?

$$\int f(g(x)) \frac{\partial g(x)}{dx} dx = \int f(u) du \qquad \text{where } u = g(x)$$

Multivariate:

$$\int f(g(\boldsymbol{x})) \left| \det \frac{\partial g(\boldsymbol{x})}{d\boldsymbol{x}} \right| d\boldsymbol{x} = \int f(\boldsymbol{u}) d\boldsymbol{u} \quad \text{where } \boldsymbol{u} = g(\boldsymbol{x})$$
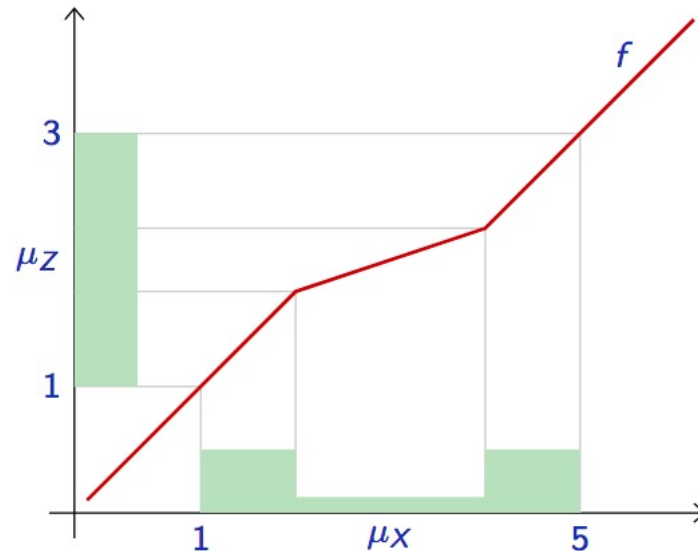
Determinant of Jacobian
of the transformation

→ Change of volume

- If $f$ is continuous, invertible, differentiable, and $x = f^{-1}(z) \equiv \phi(z)$ then

$$p_x(\boldsymbol{x}) = p_z(\boldsymbol{z}) \left| \det \left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right| \quad \text{where } \boldsymbol{x} = \phi(\boldsymbol{z})$$



The term $\left| \det \left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right|$ accounts for the local stretching of space

- If $f$ is continuous, invertible, differentiable, and $x = f^{-1}(z) \equiv \phi(z)$ then

$$p_x(\boldsymbol{x}) = p_z(\boldsymbol{z}) \left| \det \left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right| \quad \text{where } \boldsymbol{x} = \phi(\boldsymbol{z})$$

- $x = $ data we want to model, $\quad z = $ known noise

- $\phi_\theta(z)$ will be a neural network with parameters $\theta$
  - Must be continuous, invertible, differentiable

- Output of $\phi$ is a potential sample $x$
  - **Learn the right $\phi$**: adjust weights $\theta$ to maximize data probability (formula above)

# Change of Variables with Neural Networks

- If $f$ is continuous, invertible, differentiable, and $x = f^{-1}(z) \equiv \phi(z)$ then

$$p_x(\boldsymbol{x}) = p_z(\boldsymbol{z}) \left| \det\left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right| \quad \text{where } \boldsymbol{x} = \phi(\boldsymbol{z})$$

- $x$ = data we want to model,  $z$ = known noise

$\phi(\boldsymbol{z})$ neural network $\qquad\qquad\qquad \phi^{-1}(\boldsymbol{x})$ inverse

- Input  = a sample of  noise  $\iff$  – Input  = a sample X
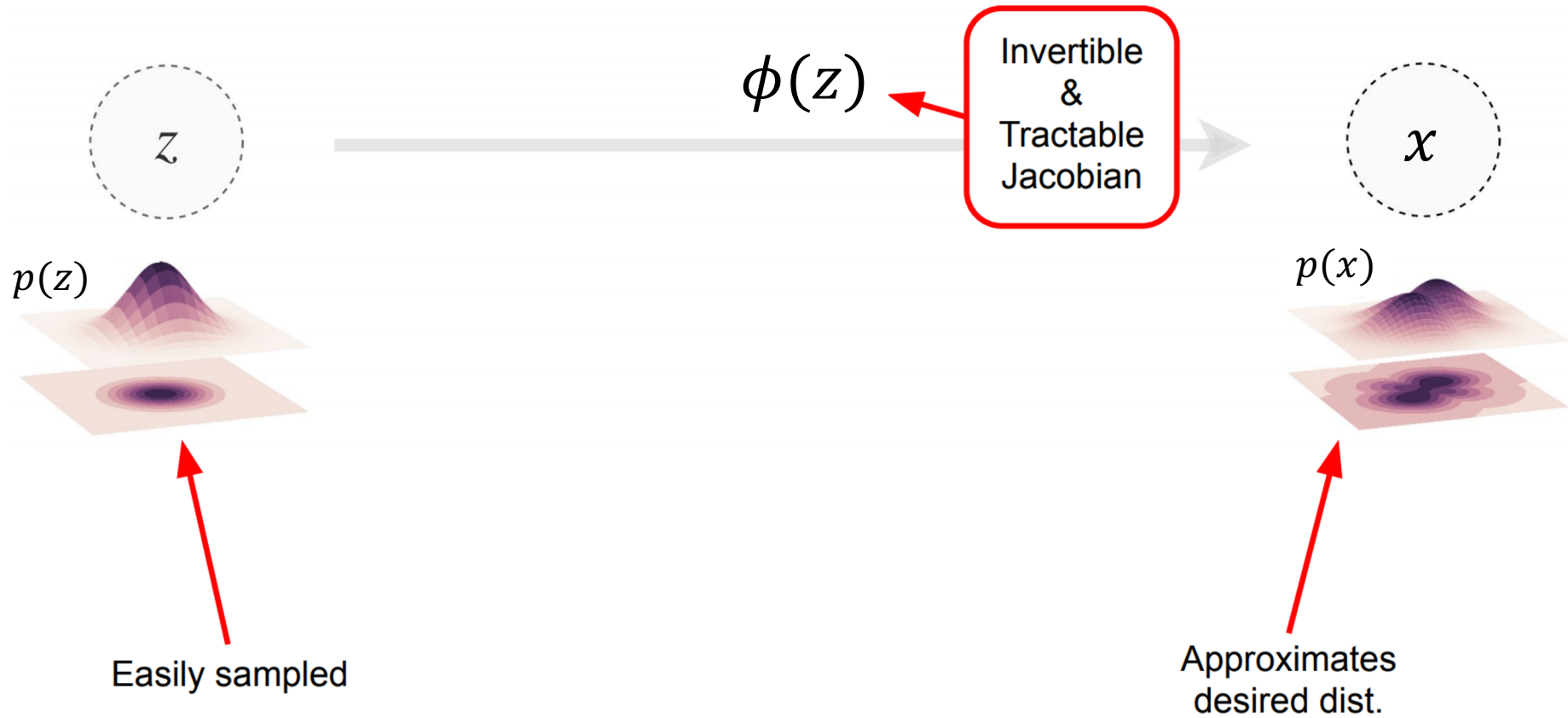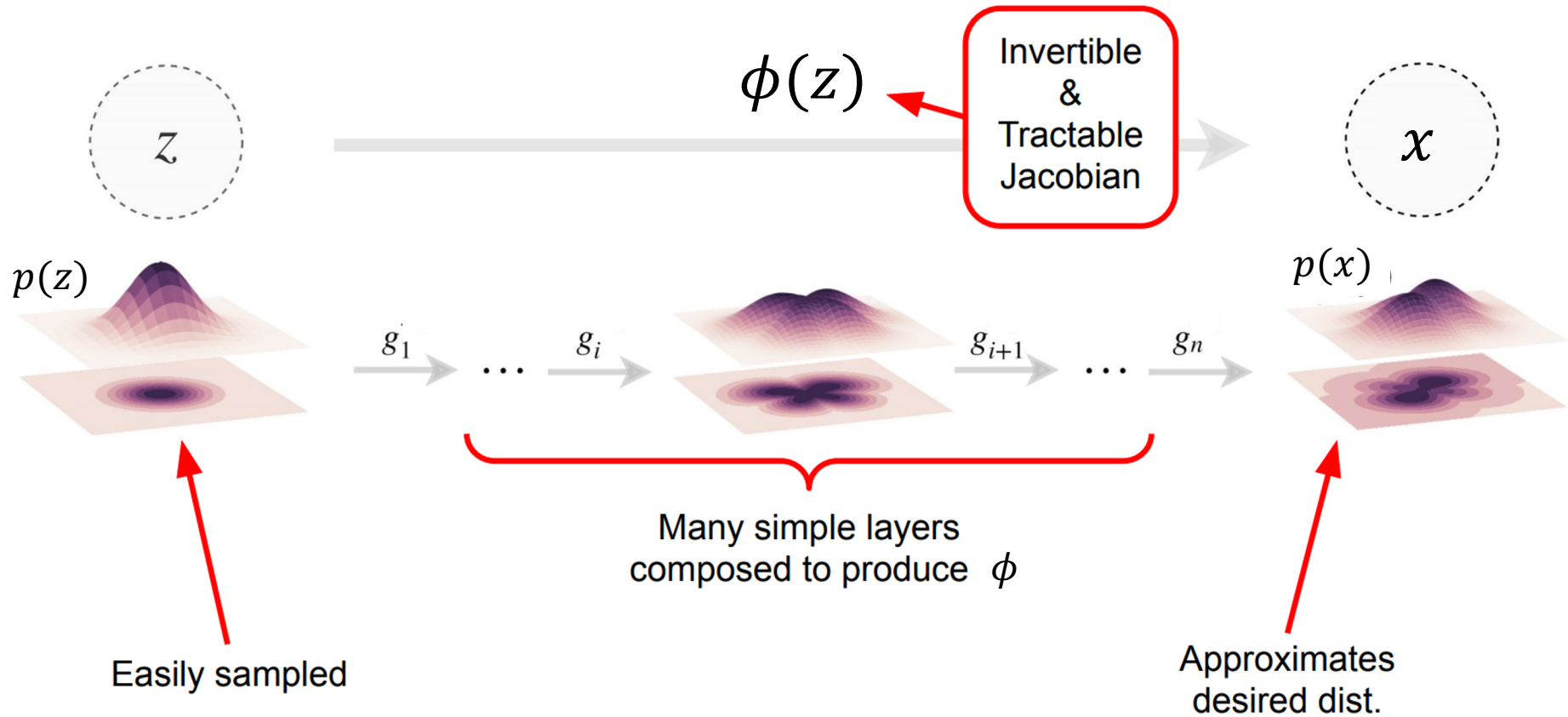- Output = a sample of  X  – Output = a sample of  noise

- Calculate the probability of  a sample using the formula above

$$p_x(\boldsymbol{x}) = p_z(\boldsymbol{z}) \left| \det\left( \frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}} \right)^{-1} \right|$$



$\phi(z)$

$z$

$x$

$p(z)$

$p(x)$

$\phi(z)$

Invertible & Tractable Jacobian

$z$

$x$

$p(z)$

$p(x)$

Easily sampled

Approximates desired dist.

$\phi(z)$

Invertible & Tractable Jacobian

$z$

$x$

$p(z)$

$p(x)$

$g_1$ ... $g_i$ $g_{i+1}$ ... $g_n$

Many simple layers composed to produce $\phi$

Easily sampled

Approximates desired dist.

# Normalizing Flows

$\phi(z)$

Invertible & Tractable Jacobian

$z$          $x$

$p(z)$          $p(x)$

$g_1$   ...   $g_i$     $g_{i+1}$   ...   $g_n$

Easily sampled      Approximates desired dist.

$$\phi^{-1}(x)$$

$$x \qquad z$$

$$\phi(z)$$

$$p_z(\phi^{-1}(x)) \left|\det\left(\frac{\partial \phi^{-1}(\boldsymbol{x})}{d\boldsymbol{x}}\right)\right|$$

$$p_x(x) \qquad p_z(z)$$

$$p_z(z) \left|\det\left(\frac{\partial \phi(\boldsymbol{z})}{d\boldsymbol{z}}\right)^{-1}\right|$$
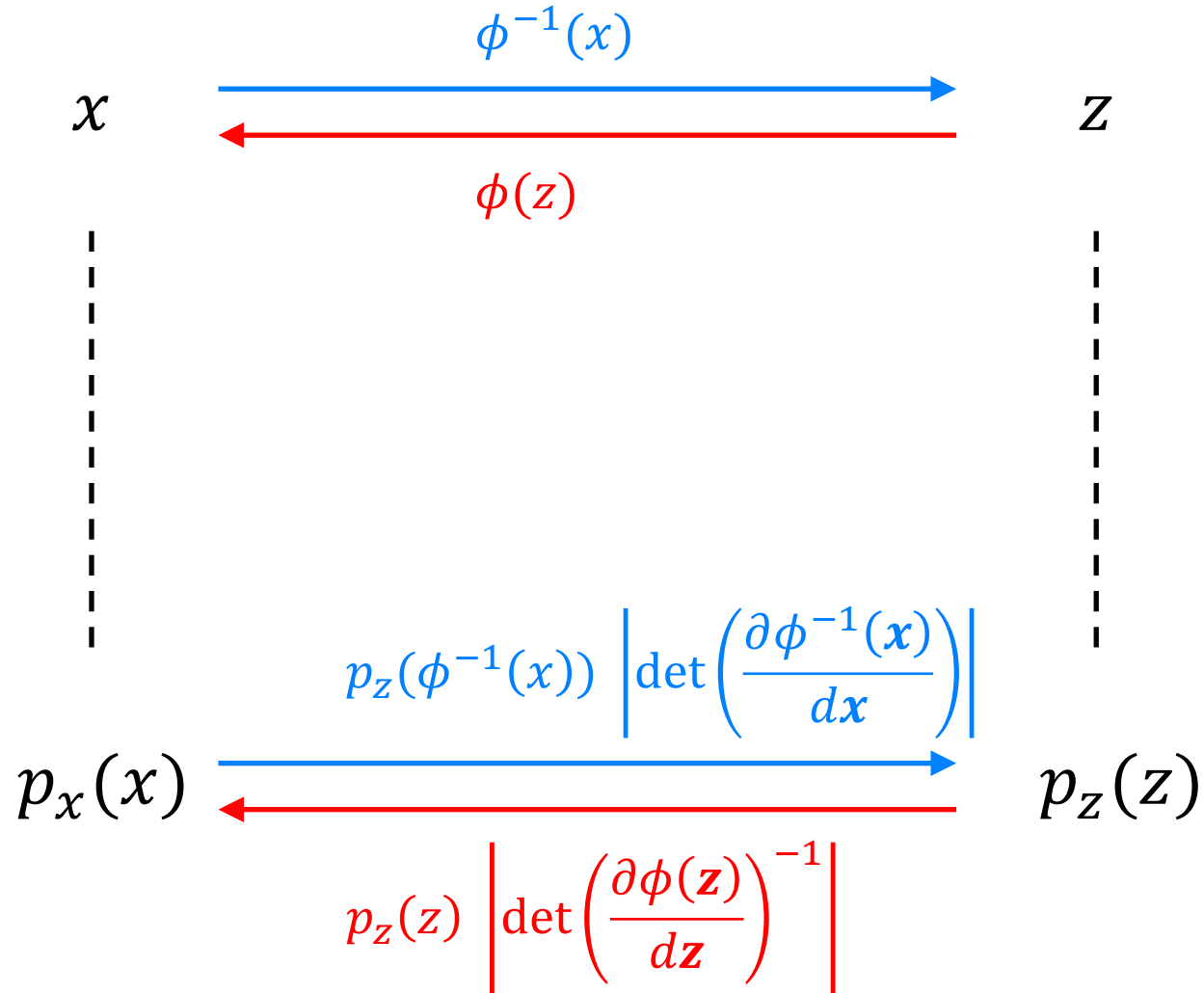
- **Learn** $\theta$ with maximum likelihood

$$\max_\theta p(x) = \max_\theta p_z\left(\phi_\theta^{-1}(x)\right) \left| \det\left( \frac{\partial \phi_\theta^{-1}(\boldsymbol{x})}{d\boldsymbol{x}} \right) \right|$$

  – Gradient descent on $\theta$
  – Find transformation s.t. data is most likely

- **Benefits** once trained
  – Can evaluate p(x) for any point X
  – Can generate "new" data points
    - Sample noise: $z \sim p(z)$
    - Transform: $\phi(z) = x$

# Example Normalizing Flow: Real NVP

- Data vector $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

- Transformation

  Functions f() and g() are neural networks

  $\phi(z)$:  $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \phi_1(z) \\ \phi_2(z) \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 * f(z_1) + g(z_1) \end{pmatrix}$
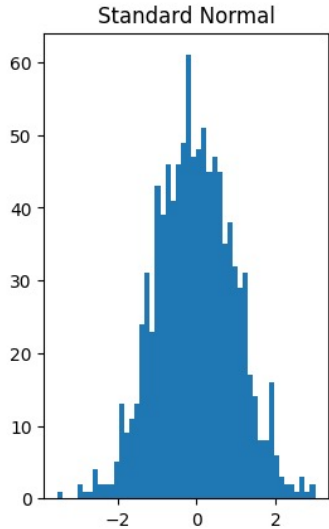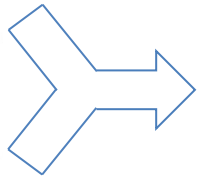
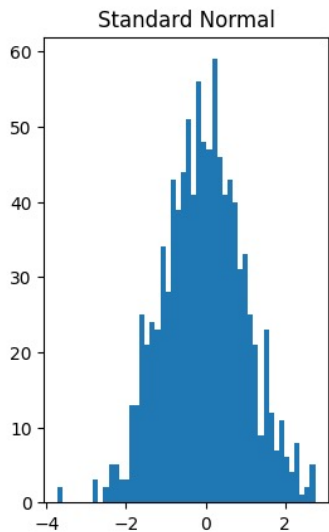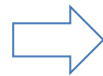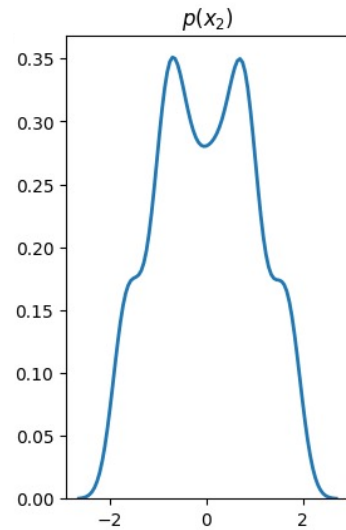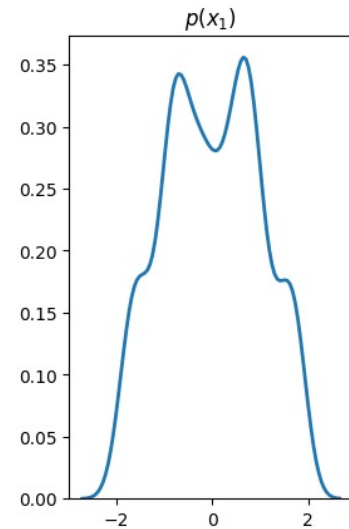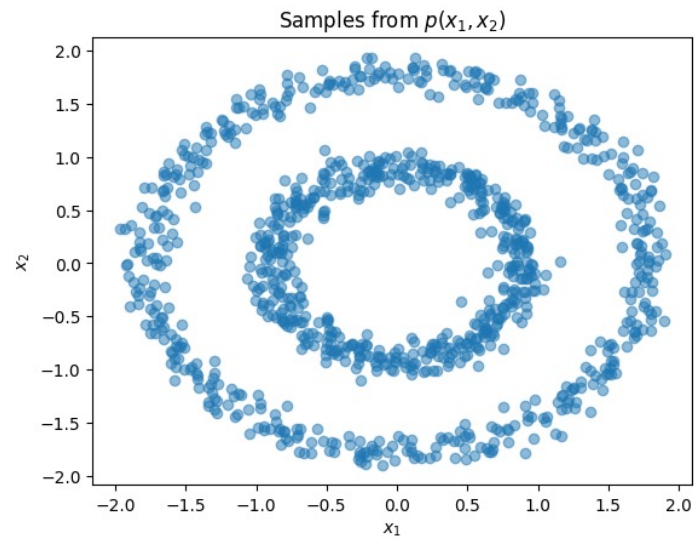  $\phi^{-1}(x)$:  $\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} \phi_1^{-1}(x) \\ \phi_2^{-1}(x) \end{pmatrix} = \begin{pmatrix} x_1 \\ (x_2 - g(x_1))/f(x_1) \end{pmatrix}$

- Determinant:

  Jacobian is lower triangular

  $$\det\left(\frac{\partial \phi(\mathbf{z})}{d\mathbf{z}}\right) = \det\left(\begin{pmatrix} 1 & 0 \\ \left(\frac{\partial \phi_2(z)}{dz_1}\right) & f(z_1) \end{pmatrix}\right) = f(z_2)$$

Standard Normal

$z_1$

$z_2$

$\phi(z)$

Samples from $p(x_1, x_2)$

$p(x_1)$

$p(x_2)$

Parameterize flow using Real NVP coupling layers

Each layer contains arbitrary neural nets

Training step

Draw samples from model

Compute loss function

Gradient descent

Desired accuracy?

Save trained model

Markov chain using samples from model

generating samples is "embarrassingly parallel"