# Variational Autoencoders (VAEs)

## DL4DS – Spring 2024

# April Dates

| Sunday | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday |
|---|---|---|---|---|---|---|
| | April 1 | 2 | 3 | 4<br>**GANs** | 5 | 6 |
| 7 | 8 | 9<br>**VAEs** | 10<br>Discussion | 11<br>**Diffusion Models** | 12 | 13 |
| 14 | 15 | 16<br>**Graph Neural Nets (VizWiz Leaders Share)** | 17<br>Discussion | 18<br>**Reinforcement Learning** | 19 | 20 |
| 21 | 22 | 23<br>**TBD/Overflow (JEPA Models)** | 24<br>Discussion | 25<br>★ **Project Presentations 1** ★ | 26 | 27 |
| 28 | 29 | 30<br>★ **Project Presentations 2** ★ | May 1<br>**Discussion??** | 2<br>Study Period | 3<br>Study Period | 4 |
| 5 | 6<br>Final Exams | 7 | 8<br>Final report & Repo ** | 9 | 10 | 11 |

** Might be earlier. Depends on when grades are due.

# Diederik P. Kingma

**Auto-Encoding Variational Bayes**

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
dpkingma@gmail.com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

2016 OpenAI, founding member
2017 PhD U. of Amsterdam
2018– Google DeepMind

## Diederik P. Kingma

Other names ▸
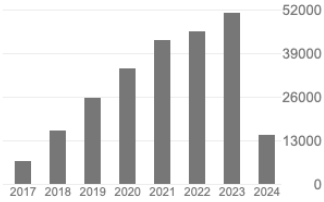
Research Scientist, Google Brain
Verified email at google.com - Homepage

Machine Learning    Deep Learning    Neural Networks
Generative Models    Artificial Intelligence

FOLLOWING

GET MY OWN PROFILE

Cited by

| | All | Since 2019 |
|---|---|---|
| Citations | 241353 | 214654 |
| h-index | 37 | 36 |
| i10-index | 39 | 39 |

| TITLE | CITED BY | YEAR |
|---|---|---|
| Adam: A method for stochastic optimization<br>DP Kingma, J Ba<br>arXiv preprint arXiv:1412.6980 | 180174 | 2014 |
| Auto-Encoding Variational Bayes<br>DP Kingma, M Welling<br>arXiv preprint arXiv:1312.6114 | 35092 | 2013 |
| Semi-Supervised Learning with Deep Generative Models<br>DP Kingma, S Mohamed, DJ Rezende, M Welling<br>Advances in Neural Information Processing Systems, 3581-3589 | 3431 | 2014 |
| Score-based generative modeling through stochastic differential equations<br>Y Song, J Sohl-Dickstein, DP Kingma, A Kumar, S Ermon, B Poole<br>arXiv preprint arXiv:2011.13456 | 3126 | 2020 |
| Glow: Generative Flow with Invertible 1x1 Convolutions<br>DP Kingma, P Dhariwal<br>Advances in Neural Information Processing Systems, 10215-10224 | 3097 | 2018 |
| An Introduction to Variational Autoencoders<br>DP Kingma, M Welling<br>Foundations and Trends® in Machine Learning 12 (4), 307-392 | 2433 | 2019 |

Public access    VIEW ALL

0 articles    3 articles
not available    available

Based on funding mandates

3

# Variational Autoencoder

Variational Inference: A method from machine learning that approximates probability densities through optimization.

Autoencoder: A type of artificial neural network used to learn efficient codings of unlabeled data in an unsupervised manner.

VAE is an autoencoder whose encodings distribution is regularized during the training to ensure that its latent space has good properties allowing us to generate new data.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Auto-Encoding Variational Bayes

Diederik P. Kingma
Machine Learning Group
Universiteit van Amsterdam
...@...com

Max Welling
Machine Learning Group
Universiteit van Amsterdam
welling.max@gmail.com

Autoencoder:  A type of artificial neural network used to learn efficient codings of unlabeled data in an unsupervised manner.

Variational Inference:  A method from machine learning that approximates probability densities through optimization.
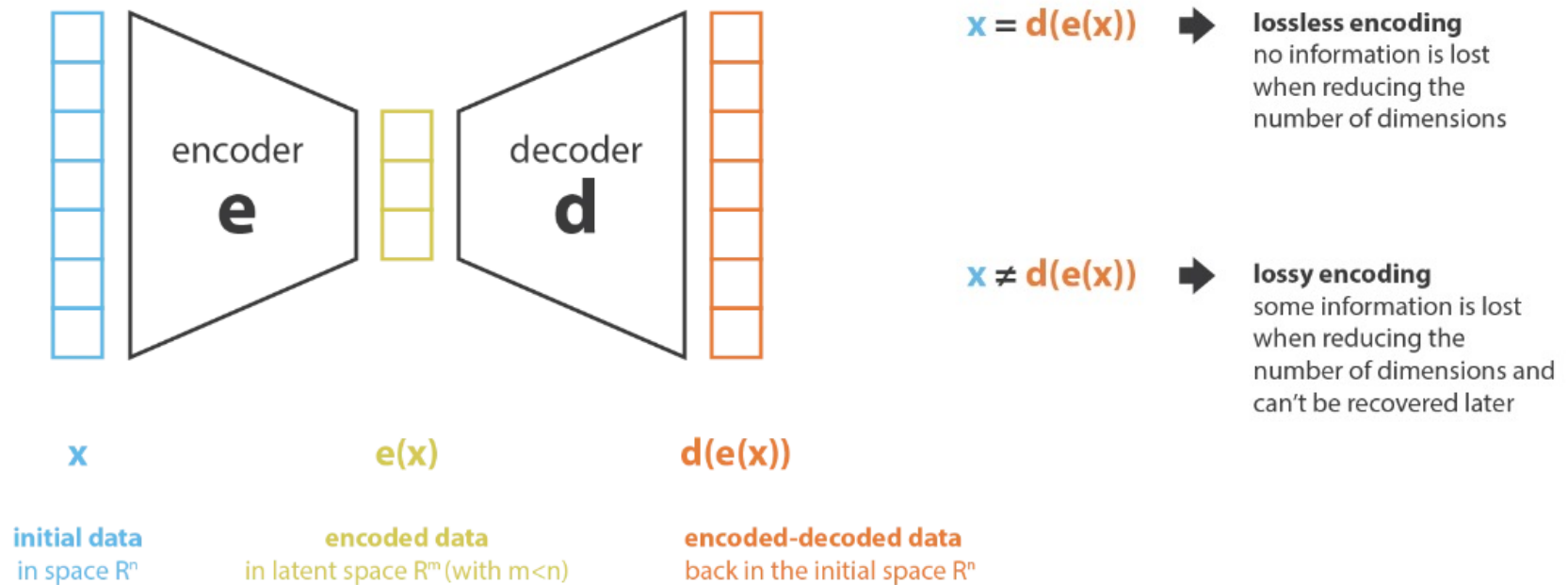
Bayesian since joint density is decomposed into prior and posterior density distributions using Bayes Rule:

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{x}|\mathbf{z})\, p(\mathbf{z})$$

5

# Outline

- <span style="color:red">Autoencoder and its limitations</span>
- Intuition behind VAEs
- Derivation of VAE
- Example applications

# Dimensionality reduction with an autoencoder



$x = d(e(x))$ ➡ **lossless encoding**
no information is lost when reducing the number of dimensions

$x \neq d(e(x))$ ➡ **lossy encoding**
some information is lost when reducing the number of dimensions and can't be recovered later

$x$

**initial data**
in space $R^n$

$e(x)$

**encoded data**
in latent space $R^m$ (with $m<n$)

$d(e(x))$

**encoded-decoded data**
back in the initial space $R^n$

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Dimensionality reduction with an autoencoder



**x**
initial data
in space $R^n$

**e(x)**
encoded data
in latent space $R^m$ (with m<n)

**d(e(x))**
encoded-decoded data
back in the initial space $R^n$

We want to find the best encoder, **e**, and decoder, **d**, to minimize the error between x and d(e(x)).
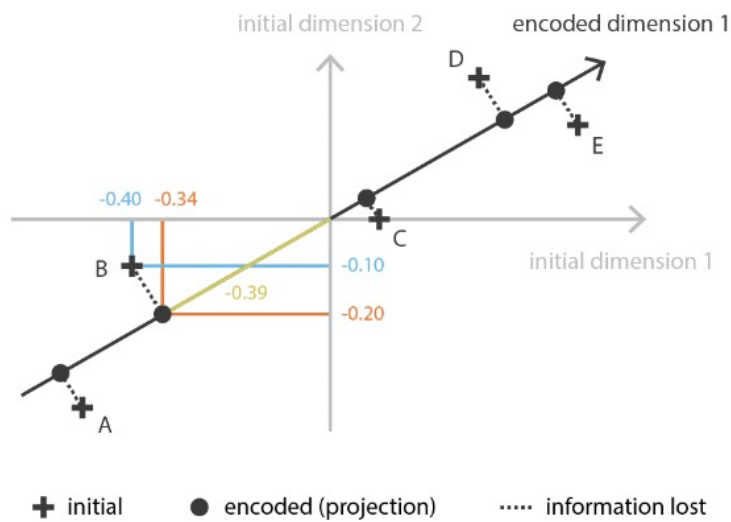
$$(e^*, d^*) = \underset{(e,d) \in E \times D}{\operatorname{argmin}} \; \epsilon(x, d(e(x)))$$

where

$$\epsilon(x, d(e(x)))$$

is the reconstruction error.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Dimensionality reduction with Principal Component Analysis (PCA)



initial dimension 2     encoded dimension 1

-0.40 -0.34

-0.10

-0.39

-0.20

initial dimension 1

+ initial    ● encoded (projection)    ⋯ information lost

$$n_d = 2 \quad n_e = 1$$

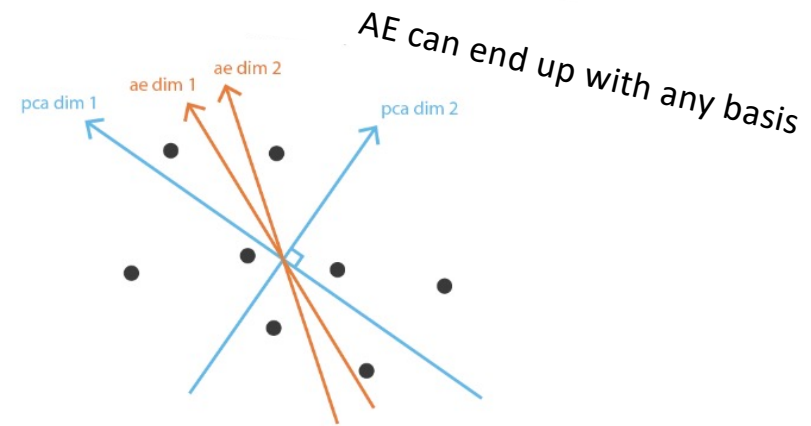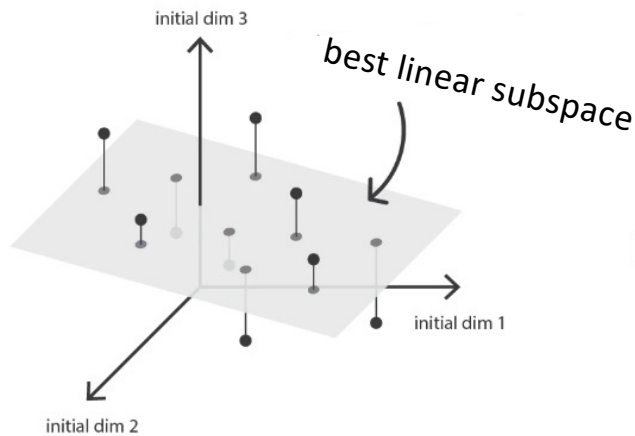| Point | Initial | Encoded | Decoded |
|-------|---------|---------|---------|
| A | (-0.50, -0.40) | -0.63 | (-0.54, -0.33) |
| B | (-0.40, -0.10) | -0.39 | (-0.34, -0.20) |
| C | (0.10, 0.00) | 0.09 | (0.07 0.04) |
| D | (0.30, 0.30) | 0.41 | (0.35, 0.21) |
| E | (0.50, 0.20) | 0.53 | (0.46, 0.27) |

Project the $n_d$-dimensional features onto an orthogonal $n_e$-dimensional subspace that minimizes Euclidean distance.

Linear Transformation!!

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Neural Network Autoencoder – 1 Linear Layer



neural network encoder

neural network decoder

x

$z = e(x)$

$\hat{x} = d(z)$

We could define encoder and decoder to each have one linear layer (no activation function), but it wouldn't necessarily converge during training to PCA solution.



initial dim 3

best linear subspace

initial dim 1

initial dim 2

AE can end up with any basis

pca dim 1

ae dim 1

ae dim 2

pca dim 2

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Neural Network Autoencoder



$$\text{loss} = \| x - \hat{x} \|^2 = \| x - d(z) \|^2 = \| x - d(e(x)) \|^2$$

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Autoencoder Reconstruction



Trained on CelebA dataset.

Kana, "Variational Autoencoders (VAEs) for Dummies -- Step by Step Tutorial", 2020

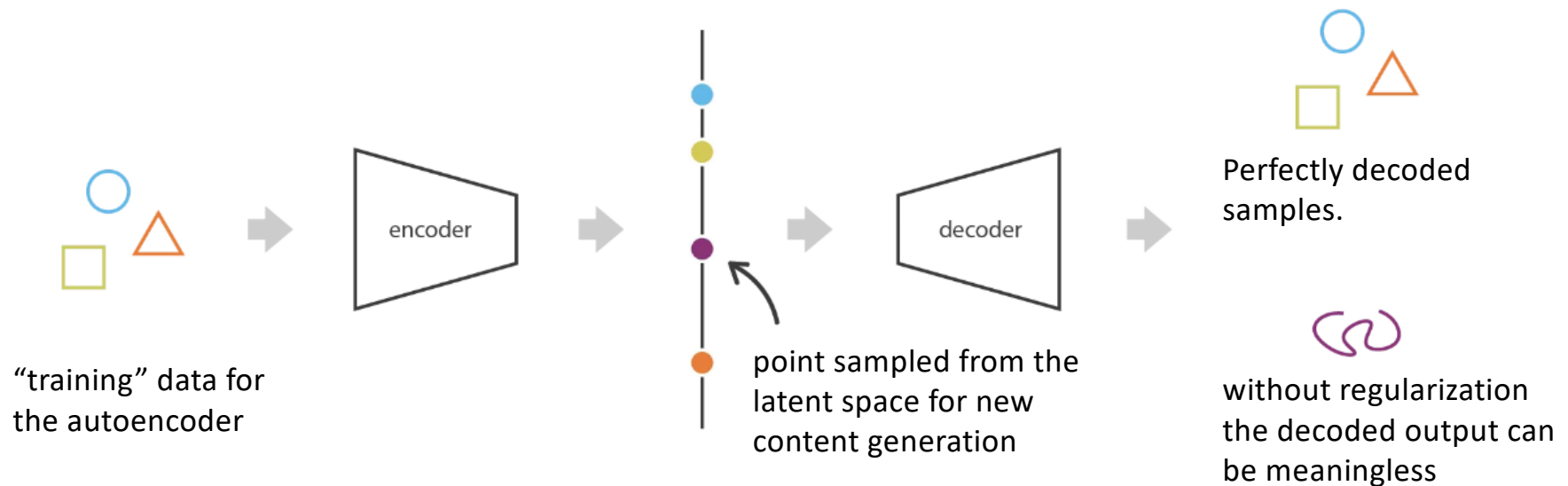# Can we generate new samples with autoencoder?



Train encoder and decoder as autoencoder.

Randomly select a different point in the latent space.

Provide as input to the decoder to generate an output.

Will this produce a good quality output? Why?

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Extreme case: Memorization



"training" data for
the autoencoder

point sampled from the
latent space for new
content generation

Perfectly decoded
samples.

without regularization
the decoded output can
be meaningless

Encoder and decoder are so powerful that they can fully
memorize the data.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019
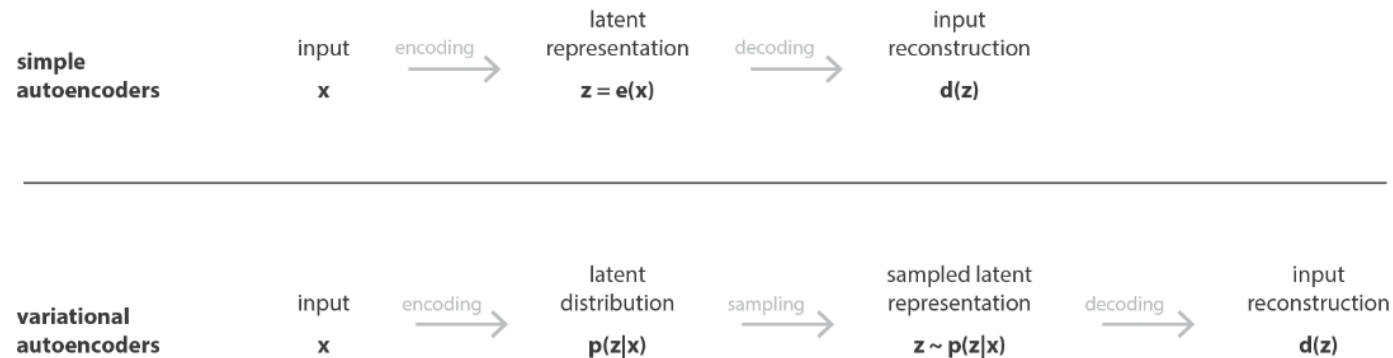
# Outline

- Autoencoder and its limitations
- <span style="color:red">Intuition behind VAEs</span>
- Derivation of VAE
- Example applications
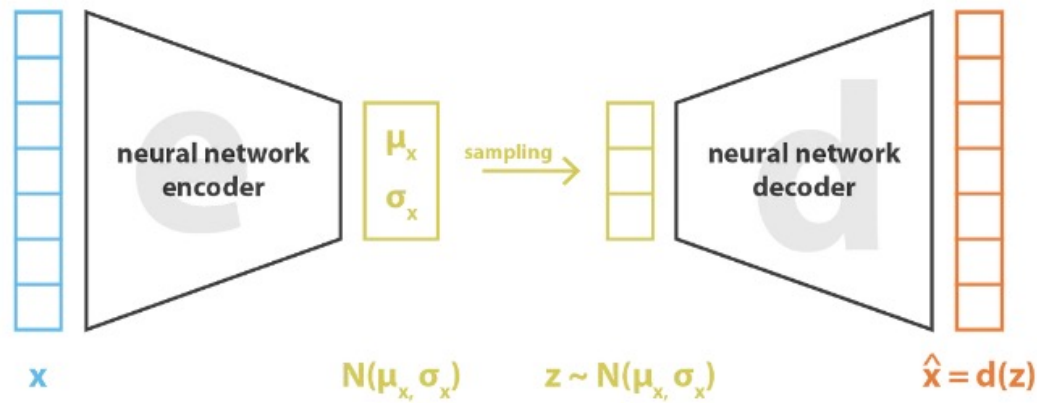
# *Variational* Autoencoder

Is an autoencoder whose training is *regularized* to avoid overfitting and ensure that the *latent space has good properties* that enable generative process.

Instead of encoding as a *single point*, encode it as a *distribution* over the latent space.

<mark>Assume distributions are normal.</mark>

| | input | encoding | latent representation | decoding | input reconstruction |
|---|---|---|---|---|---|
| **simple autoencoders** | x | → | z = e(x) | → | d(z) |

| | input | encoding | latent distribution | sampling | sampled latent representation | decoding | input reconstruction |
|---|---|---|---|---|---|---|---|
| **variational autoencoders** | x | → | p(z\|x) | → | z ~ p(z\|x) | → | d(z) |

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Autoencoder



$$\text{loss} \;=\; ||\,x - \hat{x}\,||^2 \;+\; KL[\,N(\mu_x, \sigma_x),\, N(0, I)\,] \;=\; ||\,x - d(z)\,||^2 \;+\; KL[\,N(\mu_x, \sigma_x),\, N(0, I)\,]$$

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Autoencoder



Encoder is emitting $\mu_x$ vector and $\sigma_x$ diagonal vector for independent gaussians densities.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Autoencoder



We then sample z from the multivariate Normal.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Autoencoder



Then input z to the decoder network to produce output.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Autoencoder



$$loss = \|x - \hat{x}\|^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,] = \|x - d(z)\|^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,]$$

L2 Loss    Kulback-Leibler divergence

The loss is now the L2 loss as with the autoencoder, but with an additional KL-divergence term as regularizer.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Intuitions about Regularization



point in the latent space that produce meaningless decoded output

points that are close in latent space but produce dissimilar decoded outputs

points that are close in latent space produce similar decoded outputs

irregular latent space ❌

✔ regular latent space

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Encoding to Normal distributions is not enough



without regularization ✖ ✔ with regularization

> We have to regularize the means and the covariances too!
> Regularize to a standard normal.

$$\text{loss} = \| x - \hat{x} \|^2 + \text{KL}[\, N(\mu_x, \sigma_x), N(0, I) \,]$$

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019
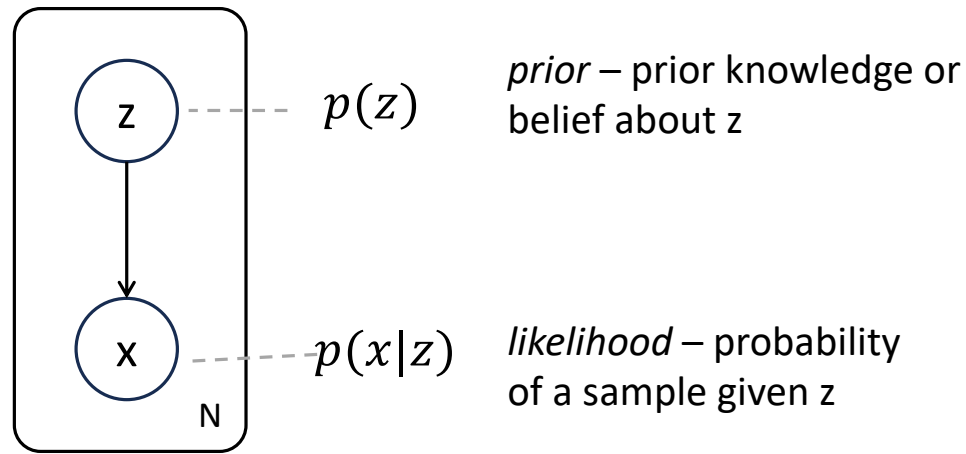
# Benefit of regularization



The continuity and completeness obtained from regularization tends to create a "gradient" over the information encoded in latent space.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

Dall-E 3

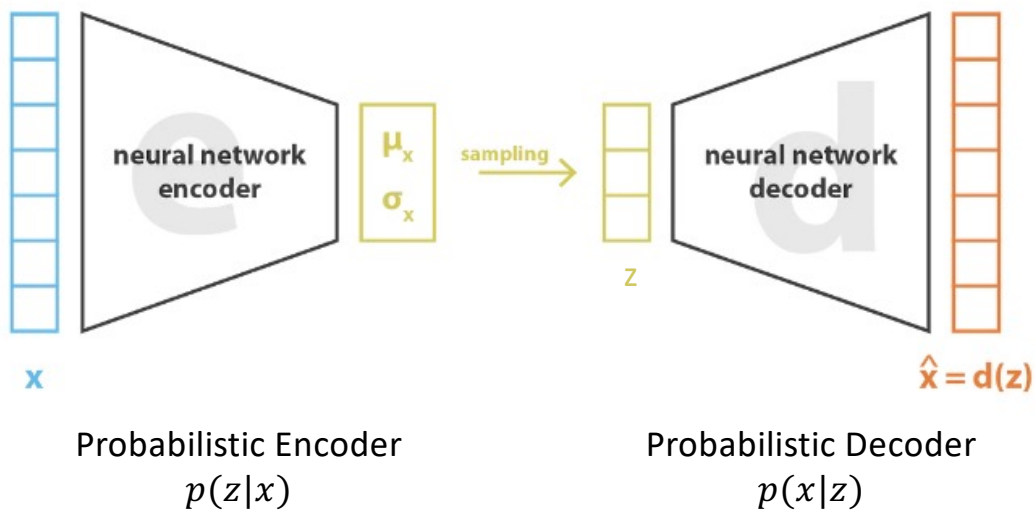# Outline

- Autoencoder and its limitations
- Intuition behind VAEs
- <span style="color:red">Derivation of VAE</span>
- Example applications

# Preliminaries: Bayesian Models

$p(z)$

*prior* – prior knowledge or belief about z

$p(x|z)$

*likelihood* – probability of a sample given z

N

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Bayesian Inference



Probabilistic Encoder
$p(z|x)$

Probabilistic Decoder
$p(x|z)$

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Bayesian Inference



Probabilistic Encoder
$p(z|x)$

Probabilistic Decoder
$p(x|z)$

*posterior* – update our knowledge of z given a new sample

*likelihood* – probability of a sample given z
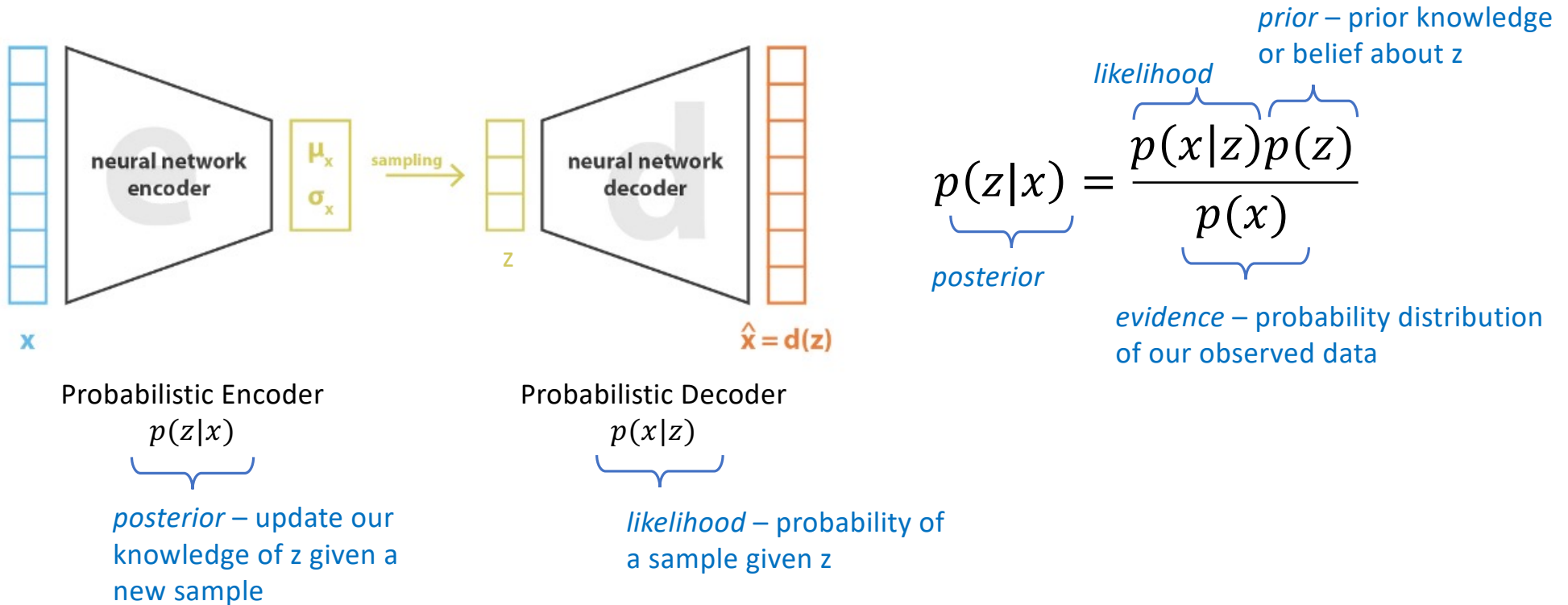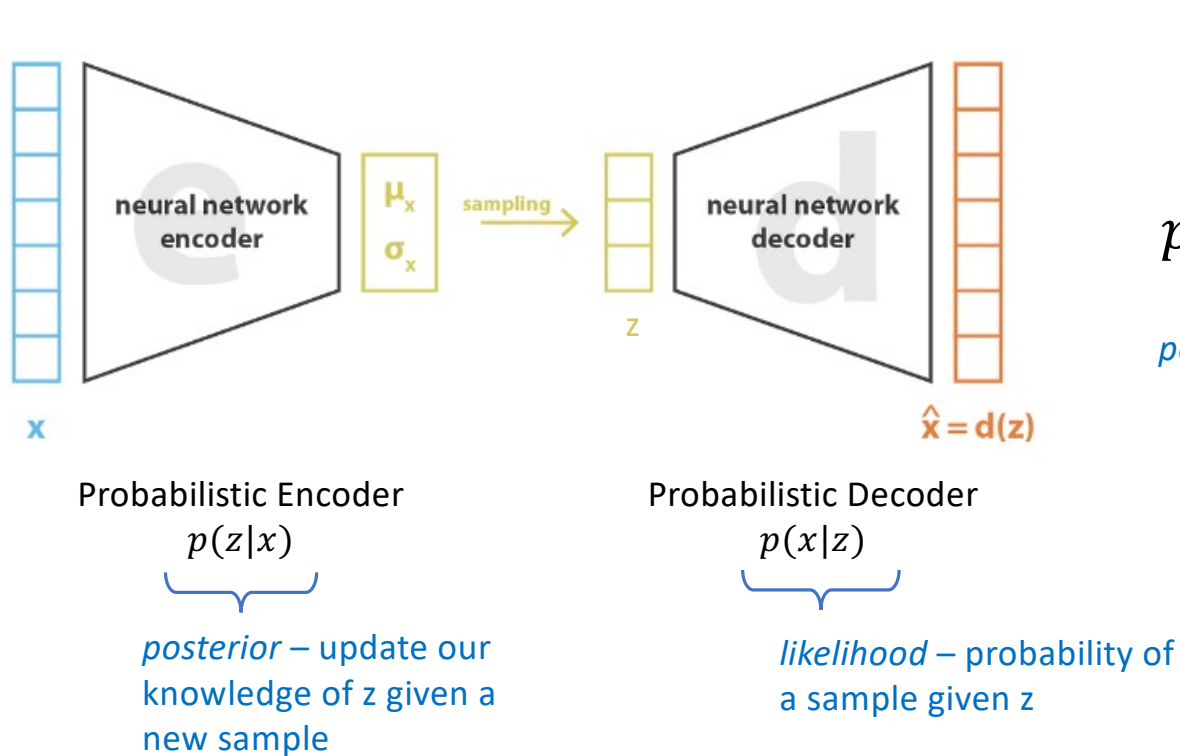
$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

We can relate the *posterior* to the *likelihood* via **Bayes Theorem.**

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Bayesian Inference



neural network encoder

$\mu_x$
$\sigma_x$

sampling

z

neural network decoder

$\hat{x} = d(z)$

Probabilistic Encoder
$p(z|x)$

*posterior* – update our knowledge of z given a new sample

Probabilistic Decoder
$p(x|z)$

*likelihood* – probability of a sample given z

*likelihood*

*prior* – prior knowledge or belief about z

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

*posterior*

*evidence* – probability distribution of our observed data

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Bayesian Inference



Probabilistic Encoder
$p(z|x)$

*posterior* – update our knowledge of z given a new sample

Probabilistic Decoder
$p(x|z)$

*likelihood* – probability of a sample given z

*likelihood*

*prior* – prior knowledge or belief about z

$$p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

*posterior*

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

We can't calculate the integral directly, but we can approximate it using *variational inference*

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Simplifying Assumptions



Probabilistic Encoder
$p(z|x)$
$\underbrace{\qquad}$
*posterior*

Probabilistic Decoder
$p(x|z)$
$\underbrace{\qquad}$
*likelihood*

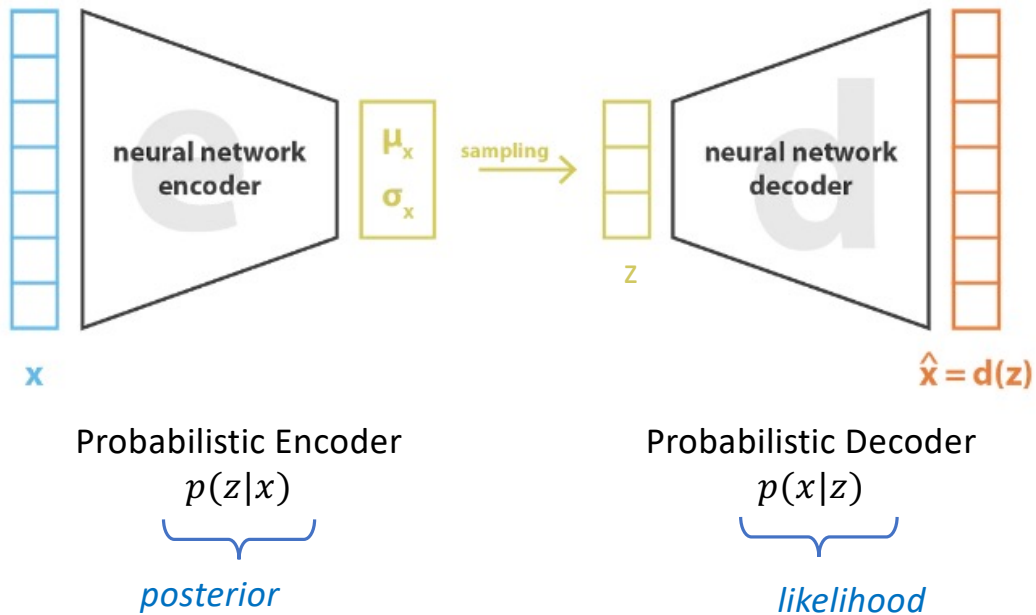Assume that the *prior* is a standard Gaussian

$$p(z) \equiv \mathcal{N}(0, I)$$

And *likelihood* is a Gaussian

$$p(x|z) \equiv \mathcal{N}(f(z), cI)$$

where $f \in F$ is a family of functions we will specify later and $c > 0$.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

# Variational Inference Formulation



Probabilistic Encoder
$p(z|x)$

posterior

Probabilistic Decoder
$p(x|z)$

likelihood

We are going to approximate *posterior* to parameterized set of Gaussians.

Approximate $p(z|x)$ by a Gaussian $q_x(z)$.

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

where $g \in G$ and $h \in H$ are a family of functions we will define shortly.

Rocca, "Understanding Variational Autoencoders (VAEs)", 2019

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

# Variational Inference

$$(g^*, h^*) = \operatorname*{arg\,min}_{(g,h) \in G \times H} KL(q_x(z), p(z|x))$$

We want to find the best functions, $g$ and $h$, to minimize the KL-divergence from the posterior $p(z|x)$.

## C.5.1 Kullback-Leibler divergence

The most common measure of distance between probability distributions $p(x)$ and $q(x)$ is the *Kullback-Leibler* or KL divergence and is defined as:

$$D_{KL}\big[p(x)||q(x)\big] = \int p(x) \log \left[\frac{p(x)}{q(x)}\right] dx. \tag{C.28}$$

# Variational Inference

$$
\begin{aligned}
(g^*, h^*) &= \operatorname*{arg\,min}_{(g,h) \in G \times H} KL(q_x(z), p(z|x)) \\
&= \operatorname*{arg\,min}_{(g,h) \in G \times H} \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}\left( \log \frac{p(x|z)p(z)}{p(x)} \right) \right)
\end{aligned}
$$

➢ Rewriting KL divergence as Expectation,
➢ log of division is difference of the logs
➢ substituting for the posterior using Bayes Theorem

# Variational Inference

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

$$
\begin{aligned}
(g^*, h^*) &= \underset{(g,h) \in G \times H}{\arg\min} \ KL(q_x(z), p(z|x)) \\
&= \underset{(g,h) \in G \times H}{\arg\min} \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}\left( \log \frac{p(x|z)p(z)}{p(x)} \right) \right) \\
&= \underset{(g,h) \in G \times H}{\arg\min} \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}(\log p(z)) - \mathbb{E}_{z \sim q_x}(\log p(x|z)) + \mathbb{E}_{z \sim q_x}(\log p(x)) \right)
\end{aligned}
$$

- ➢ log of product becomes sum of logs
- ➢ log of division becomes difference of logs

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

# Variational Inference

$$
\begin{aligned}
(g^*, h^*) = \ & \underset{(g,h) \in G \times H}{\arg\min} \ KL(q_x(z), p(z|x)) \\
= \ & \underset{(g,h) \in G \times H}{\arg\min} \ \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x} \left( \log \frac{p(x|z)p(z)}{p(x)} \right) \right) \\
= \ & \underset{(g,h) \in G \times H}{\arg\min} \ (\mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}(\log p(z)) - \mathbb{E}_{z \sim q_x}(\log p(x|z)) + \mathbb{E}_{z \sim q_x}(\log p(x))) \\
= \ & \underset{(g,h) \in G \times H}{\arg\max} \ (\mathbb{E}_{z \sim q_x}(\log p(x|z)) - KL(q_x(z), p(z)))
\end{aligned}
$$

➢ negating and converting from argmin to argmax
➢ collecting terms to form KL divergence

# Variational Inference

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

$$
\begin{aligned}
(g^*, h^*) &= \underset{(g,h) \in G \times H}{\arg\min} \; KL(q_x(z), p(z|x)) \\[2mm]
&= \underset{(g,h) \in G \times H}{\arg\min} \; \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}\left( \log \frac{p(x|z)p(z)}{p(x)} \right) \right) \\[2mm]
&= \underset{(g,h) \in G \times H}{\arg\min} \; \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}(\log p(z)) - \mathbb{E}_{z \sim q_x}(\log p(x|z)) + \mathbb{E}_{z \sim q_x}(\log p(x)) \right) \\[2mm]
&= \underset{(g,h) \in G \times H}{\arg\max} \; \left( \underbrace{\mathbb{E}_{z \sim q_x}(\log p(x|z))} - \underbrace{KL(q_x(z), p(z))} \right)
\end{aligned}
$$

Maximize the expected log likelihood.

Minimize the difference between the approximate posterior and the prior.

$$q_x(z) \equiv \mathcal{N}(g(x), h(x))$$

# Variational Inference

$$
\begin{aligned}
(g^*, h^*) =\ & \underset{(g,h) \in G \times H}{\arg\min}\ KL(q_x(z), p(z|x)) \\[2mm]
=\ & \underset{(g,h) \in G \times H}{\arg\min}\ \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}\left( \log \frac{p(x|z)p(z)}{p(x)} \right) \right) \\[2mm]
=\ & \underset{(g,h) \in G \times H}{\arg\min}\ \left( \mathbb{E}_{z \sim q_x}(\log q_x(z)) - \mathbb{E}_{z \sim q_x}(\log p(z)) - \mathbb{E}_{z \sim q_x}(\log p(x|z)) + \mathbb{E}_{z \sim q_x}(\log p(x)) \right) \\[2mm]
=\ & \underset{(g,h) \in G \times H}{\arg\max}\ \left( \mathbb{E}_{z \sim q_x}(\log p(x|z)) - KL(q_x(z), p(z)) \right) \\[2mm]
=\ & \underset{(g,h) \in G \times H}{\arg\max}\ \left( \mathbb{E}_{z \sim q_x}\left( -\frac{||x - f(z)||^2}{2c} \right) - KL(q_x(z), p(z)) \right)
\end{aligned}
$$

Log of the Gaussian likelihood $p(x|z) \equiv \mathcal{N}(f(z), cI)$.
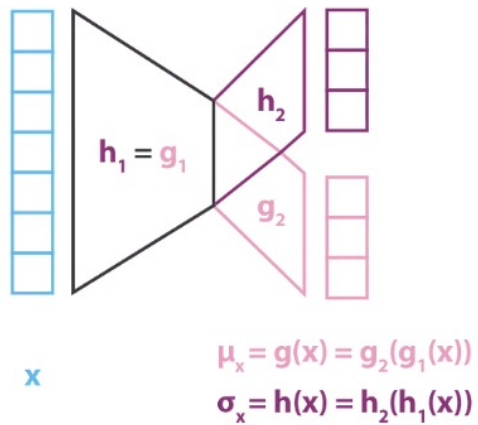
This brings our function, $f$, into the equation, so...

# Variational Inference

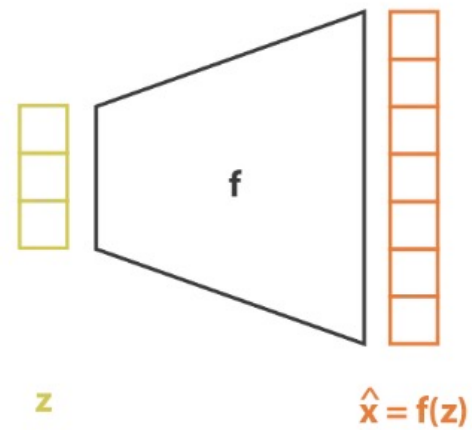We are looking for optimal $f^*$, $g^*$ and $h^*$ such that

$$(f^*, g^*, h^*) = \underset{(f,g,h) \in F \times G \times H}{\arg\max} \left( \mathbb{E}_{z \sim q_x} \left( -\frac{||x - f(z)||^2}{2c} \right) - KL(q_x(z), p(z)) \right)$$

Note that the constant, c, determines the balance between reconstruction error and the regularization term given by KL divergence.
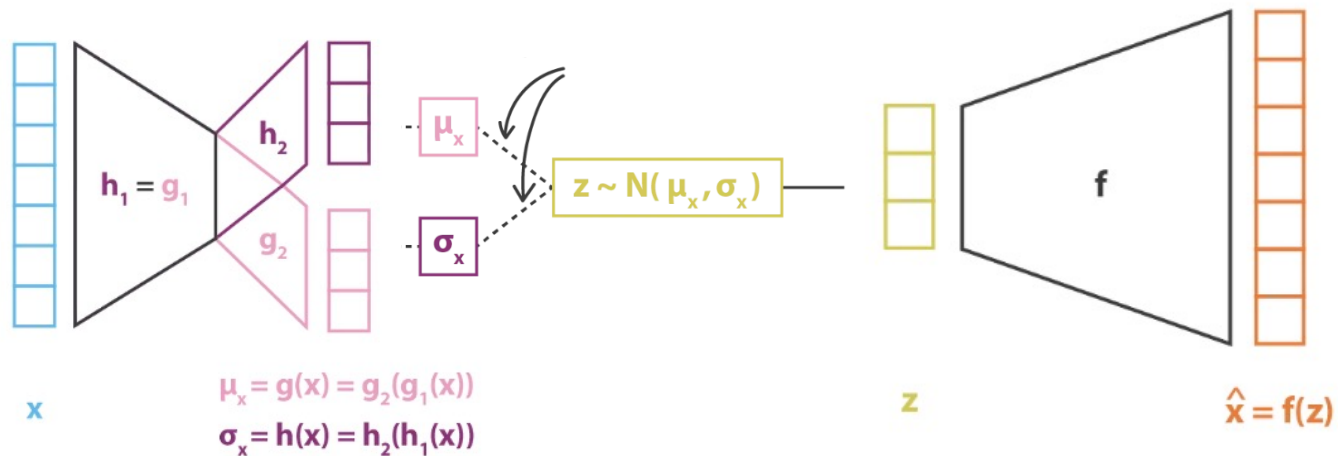
# Enter the Neural Networks



$\mu_x = g(x) = g_2(g_1(x))$

$\sigma_x = h(x) = h_2(h_1(x))$

$x$

Encoder produces the mean and variance.

$z$

$\hat{x} = f(z)$

Decoder reconstructs the input (during training)

# But one more problem to solve



$$\mu_x = g(x) = g_2(g_1(x))$$
$$\sigma_x = h(x) = h_2(h_1(x))$$

We can't backpropagate through the sampling step.

# Use the reparameterization trick



sampling prevents backpropagation and then training

$\mu_x$

$z \sim N(\mu_x, \sigma_x)$

$\sigma_x$

$\zeta \sim N(0, I)$

no backpropagation is required

$\mu_x$

$z = \sigma_x \zeta + \mu_x$

$\sigma_x$

# Putting it all together



N(0, I)

h

g

$\mu_x = g(x)$
$\sigma_x = h(x)$
$\zeta \sim N(0, I)$

x

$z = \sigma_x \zeta + \mu_x$

$\hat{x} = f(z)$

f

We use a Monte-Carlo approximation to the expectation of reconstruction loss

Convert C = 1/(2c).

$$loss = C \parallel x - \hat{x} \parallel^2 + KL[\, N(\mu_x, \sigma_x), N(0, I)\,] = C \parallel x - f(z) \parallel^2 + KL[\, N(g(x), h(x)), N(0, I)\,]$$

We have as trainable neural network!

# Probability Distribution Divergence Measures

## C.5.1 Kullback-Leibler divergence

The most common measure of distance between probability distributions $p(x)$ and $q(x)$ is the *Kullback-Leibler* or KL divergence and is defined as:

$$D_{KL}\big[p(x)||q(x)\big] = \int p(x)\log\left[\frac{p(x)}{q(x)}\right]dx. \qquad (C.28)$$

## C.5.2 Jensen-Shannon divergence

The KL divergence is not symmetric (i.e., $D_{KL}[p(x)||q(x)] \neq D_{KL}[q(x)||p(x)]$). The Jensen-Shannon divergence is a measure of distance that is symmetric by construction:

$$D_{JS}\big[p(x)||q(x)\big] = \frac{1}{2}D_{KL}\left[p(x)\Big|\Big|\frac{p(x)+q(x)}{2}\right] + \frac{1}{2}D_{KL}\left[q(x)\Big|\Big|\frac{p(x)+q(x)}{2}\right]. \qquad (C.30)$$

It is the mean divergence of $p(x)$ and $q(x)$ to the average of the two distributions.

Prince, *Understanding Deep Learning*

# Outline

- Autoencoder and its limitations
- Intuition behind VAEs
- Derivation of VAE
- Example applications

# Generating high quality images



Vahdat & Kautz (2020) "NVAE: A deep hierarchical variational autoencoder"
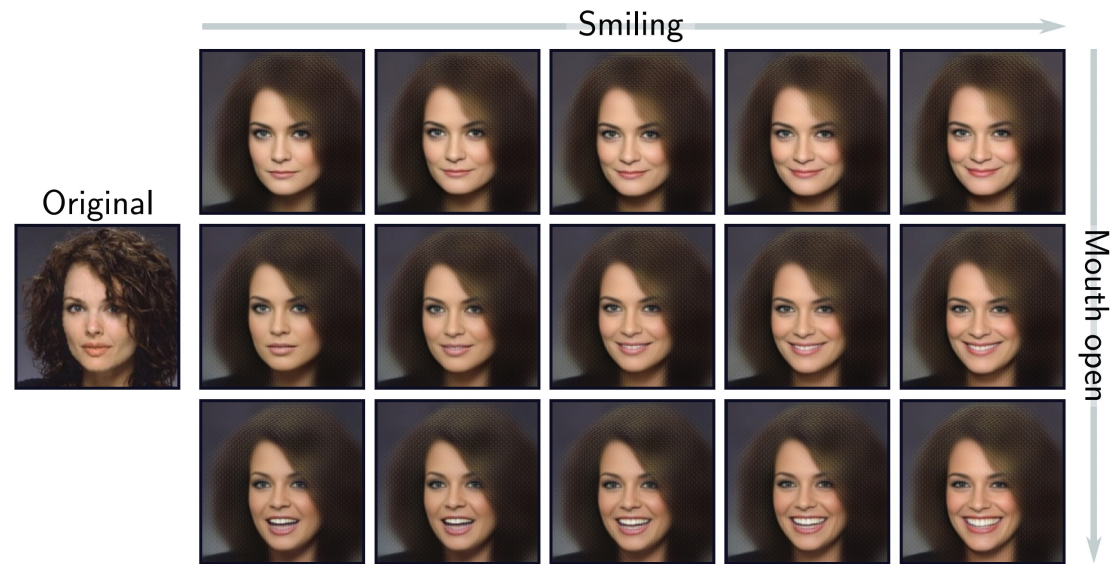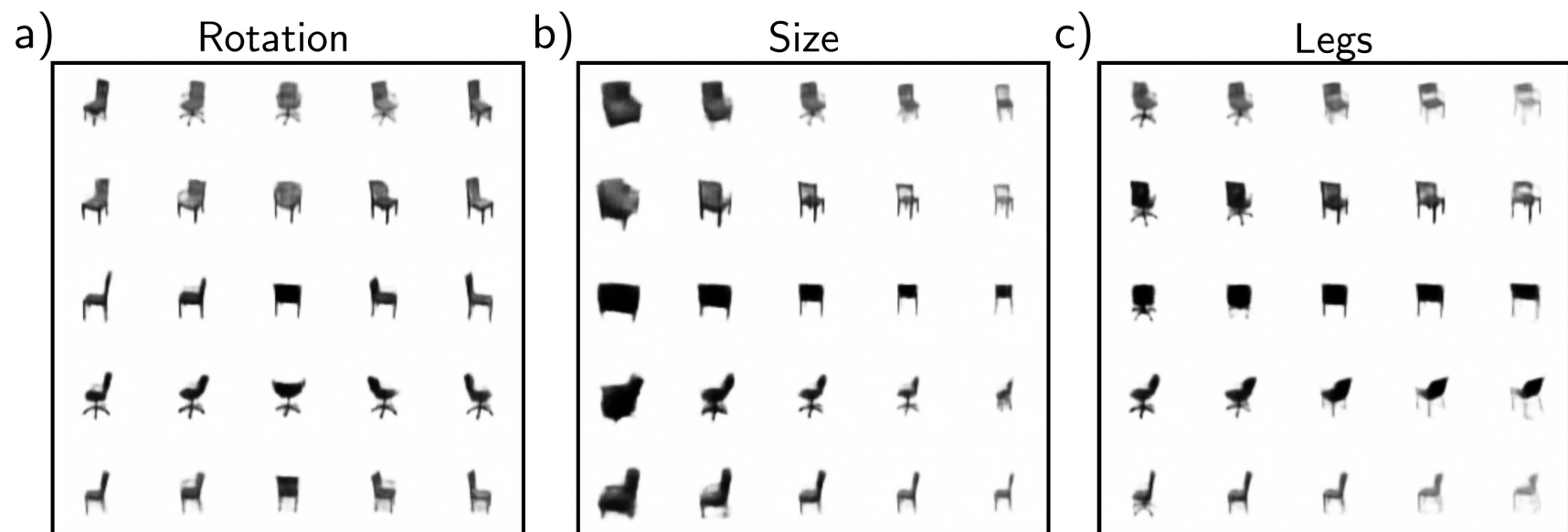
# Resynthesizing real data with changes



**Figure 17.13** Resynthesis. The original image on the left is projected into the latent space using the encoder, and the mean of the predicted Gaussian is chosen to represent the image. The center-left image in the grid is the reconstruction of the input. The other images are reconstructions after manipulating the latent space in directions representing smiling/neutral (horizontal) and mouth open/closed (vertical). Adapted from White (2016).

# Disentanglement of the latent space



a) Rotation  b) Size  c) Legs

Chen et al (2021) "Cross-layer distillation with semantic calibration."

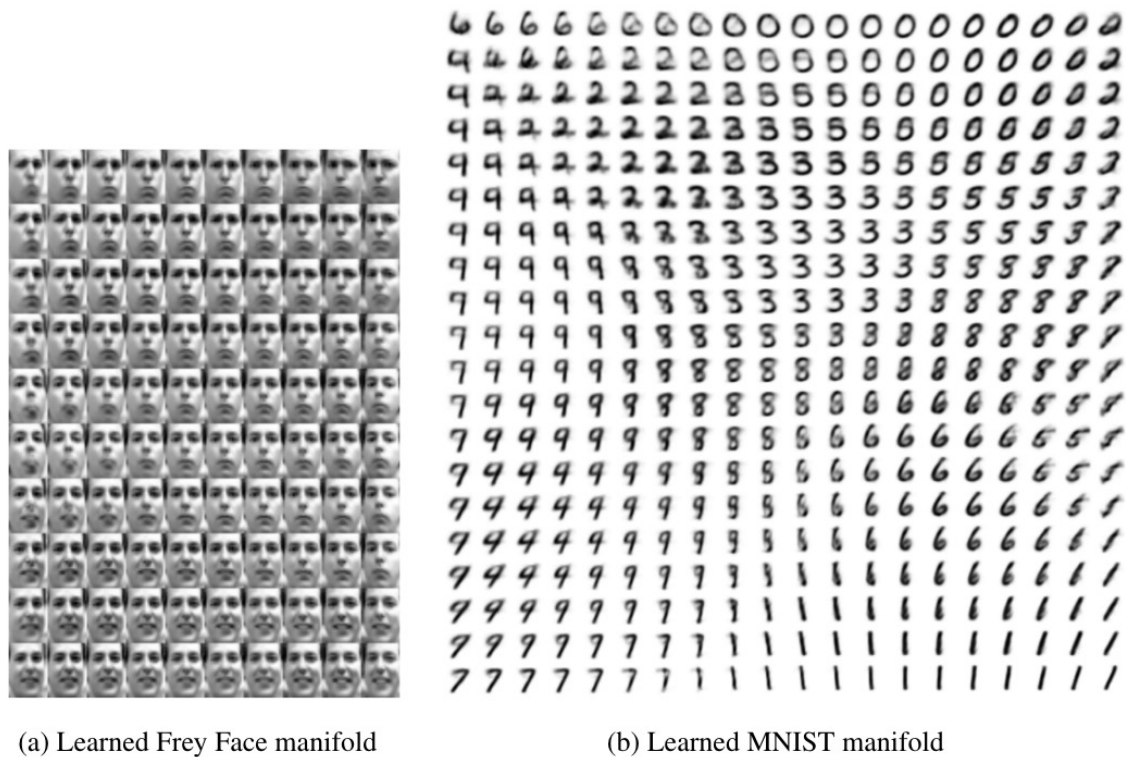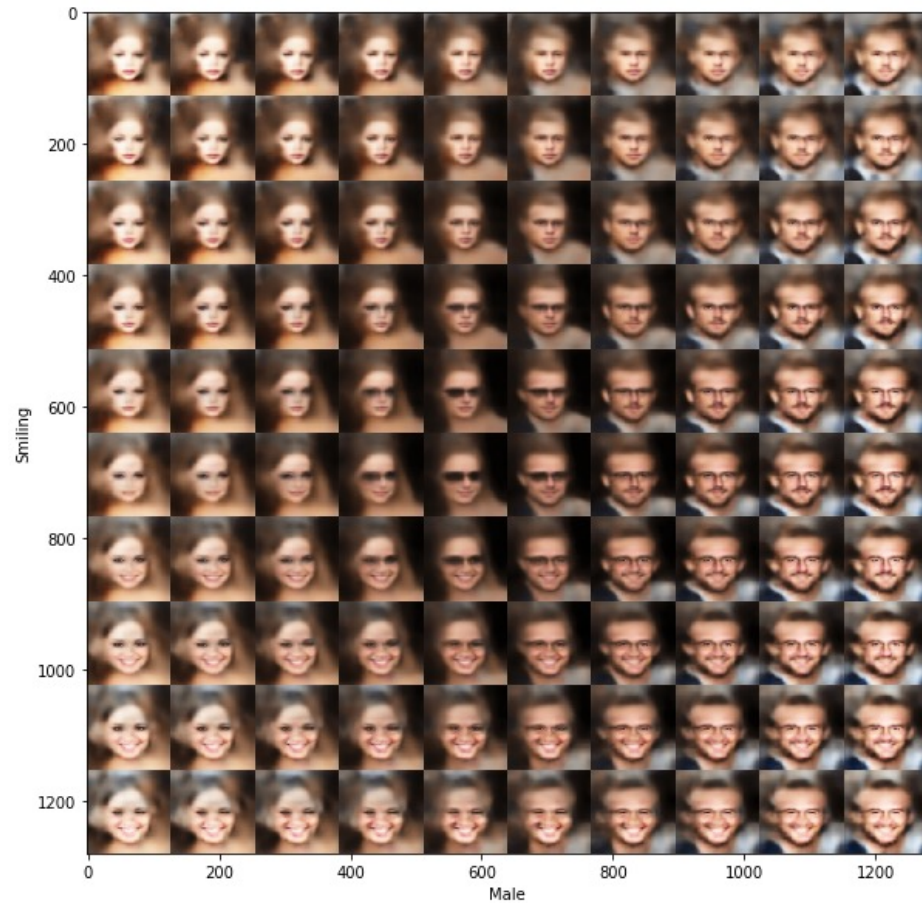(a) Learned Frey Face manifold      (b) Learned MNIST manifold

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 $p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})$ with the learned parameters $\boldsymbol{\theta}$.

# Conditional VAEs

# Debiasing



Capable of uncovering **underlying features** in a dataset

Homogeneous skin color, pose

VS

Diverse skin color, pose, illumination

How can we use this information to create fair and representative datasets?

Amini et al, "Uncovering and Mitigating Algorithmic Bias through Learned Latent Structure," 2019
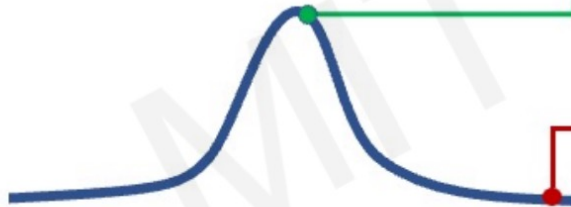
# Outlier Detection

- **Problem:** How can we detect when we encounter something new or rare?
- **Strategy:** Leverage generative models, detect outliers in the distribution
- Use outliers during training to improve even more!

**95% of Driving Data:**
(1) sunny, (2) highway, (3) straight road



Detect outliers to avoid unpredictable behavior when training
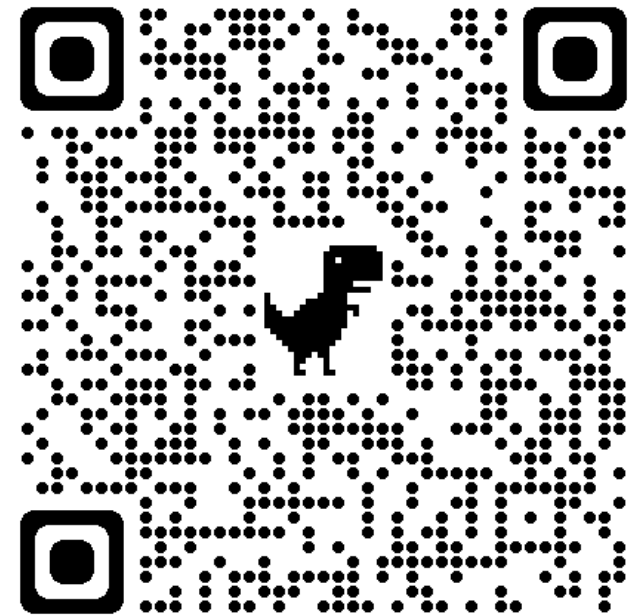
Edge Cases    Harsh Weather    Pedestrians

A. Amini et al, "Variational Autoencoder for End-to-End Control of Autonomous Driving with Novelty Detection and Training De-biasing," *2018*

# Upcoming Topics

- Diffusion Models
- Graph Neural Networks
- Reinforcement Learning

# Feedback



[Link](Link)