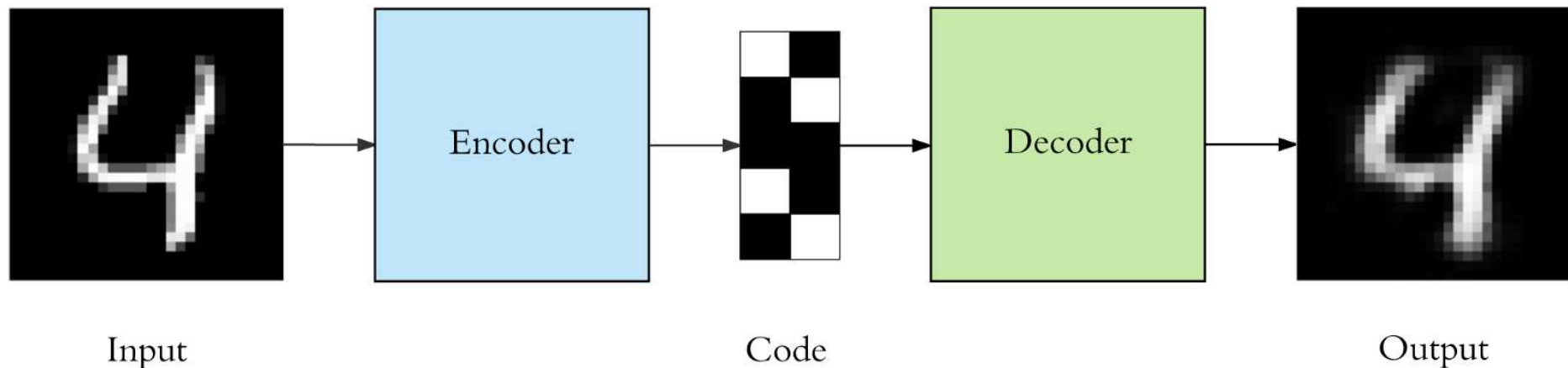# Deep Learning (CS324)

# 8. Variational autoencoders*

Prof Jianguo Zhang

SUSTech

# Generation with autoencoders?

- Can we use an autoencoder to generate data given a code?



Input    Encoder    Code    Decoder    Output
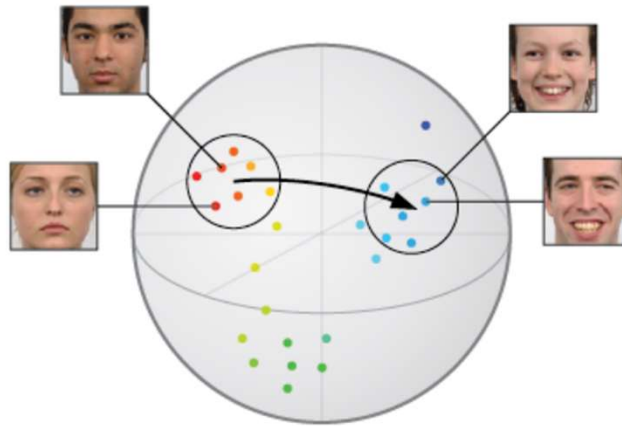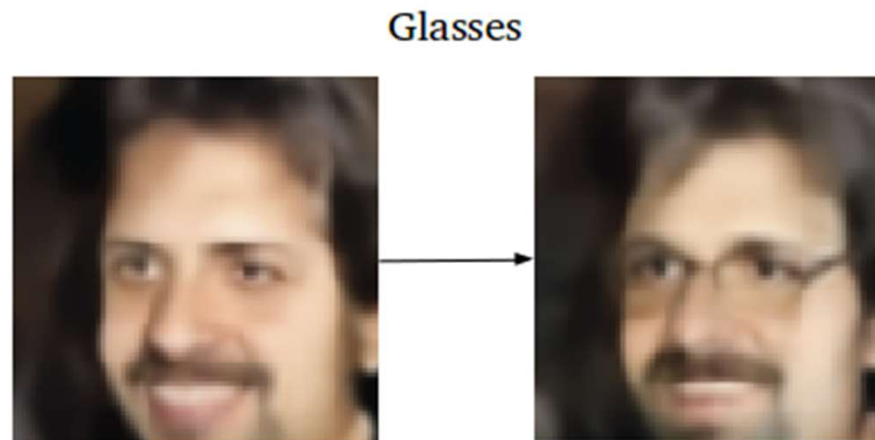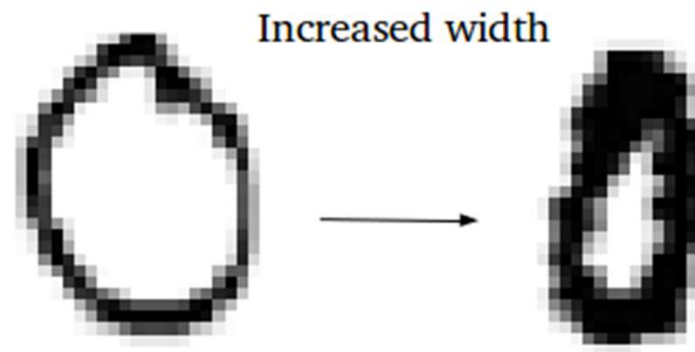
# Ok. But why? Sample application



Figure 1: Schematic of the latent space of a generative model. In the general case, a generative model includes an encoder to map from the feature space (here images of faces) into a high dimensional latent space. Vector space arithmetic can be used in the latent space to perform semantic operations. The model also includes a decoder to map from the latent space back into the feature space, where the semantic operations can be observed. If the latent space transformation is the identity function we refer to the encoding and decoding as a reconstruction of the input through the model.
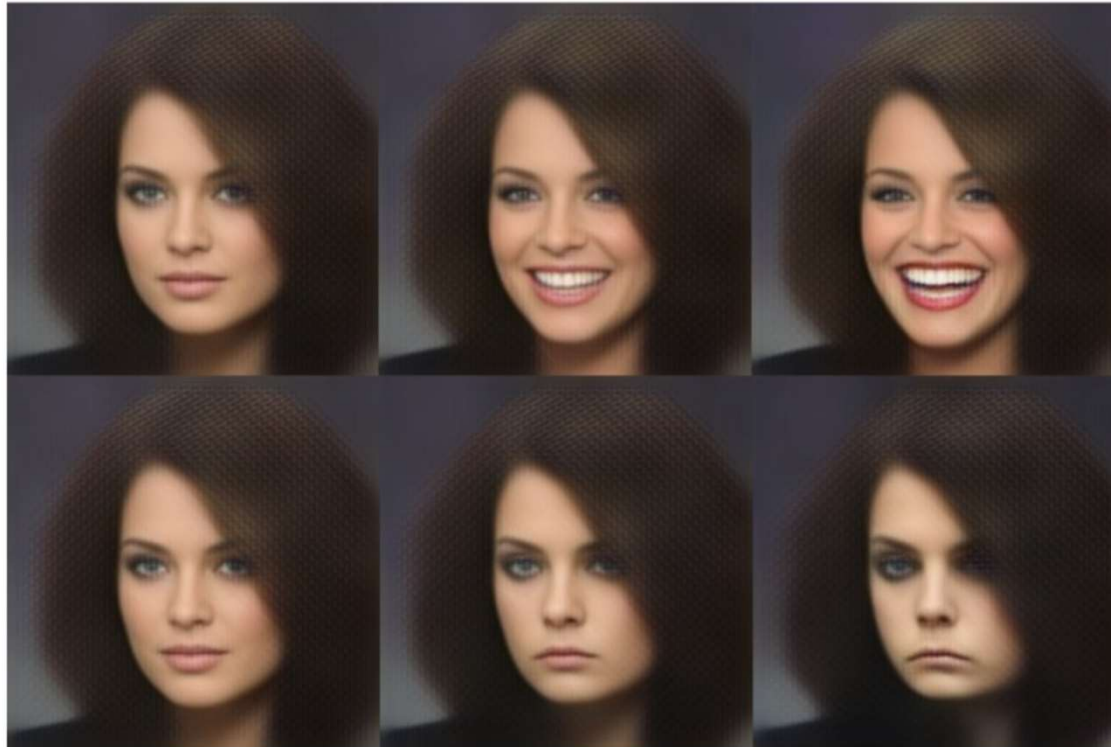
Exploring a specific variation of input data

# Ok. But why? Sample application



Increased width

Glasses

Exploring a specific variation of input data
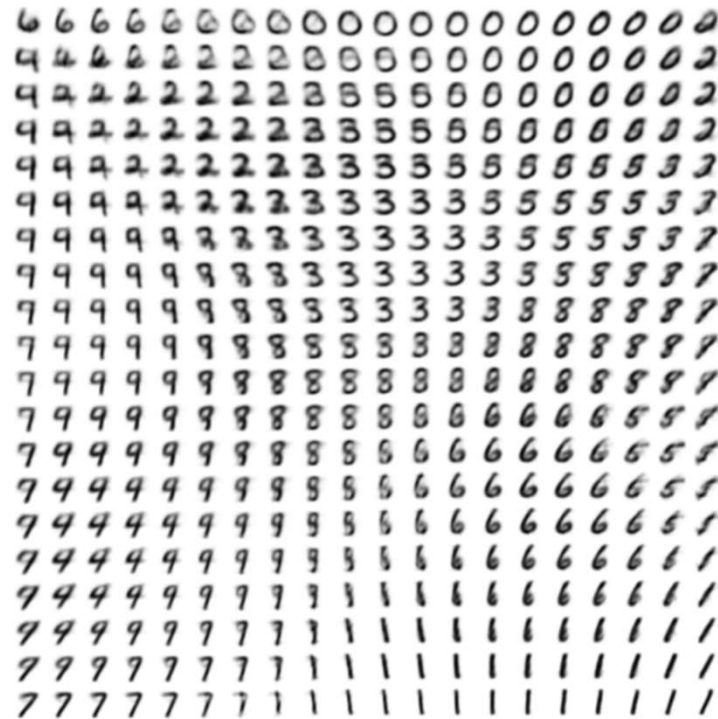
# Ok. But why? Sample application



Exploring a specific variation of input data

# Ok. But why? Sample application
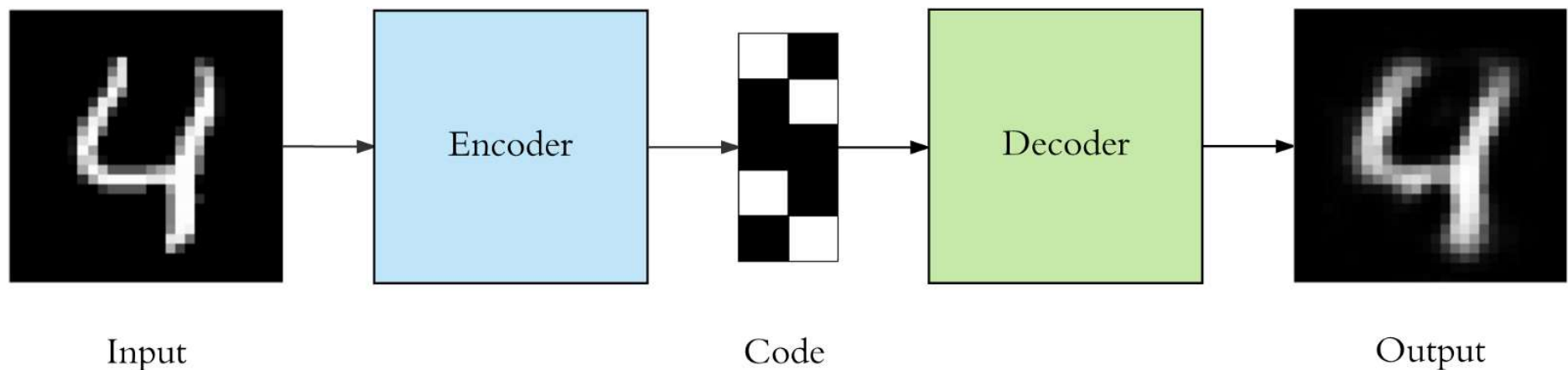


(a) Learned Frey Face mani-
fold

(b) Learned MNIST manifold

Figure 2.7: Visualizations 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 val-
ues 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 $\boldsymbol{\theta}$.
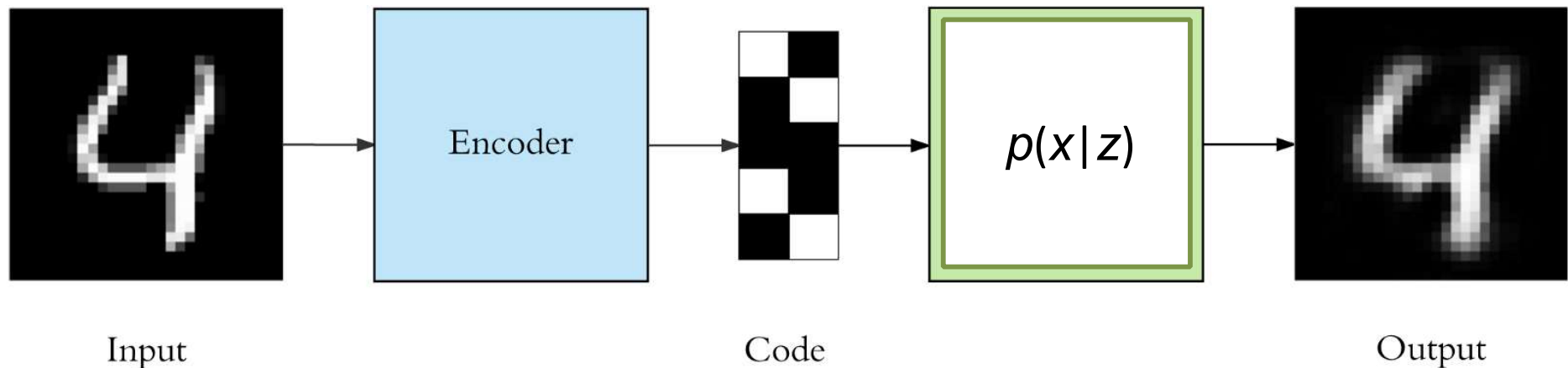
# Generation with autoencoders?

- Can we use an autoencoder to generate data given a code?



Kingma and Welling, Auto-encoding Variational Bayes, NIPS, 2013

# Generation with autoencoders?

- Can we use an autoencoder to generate data given a code?

- Yes, but…



Input        Encoder    Code    $p(x|z)$    Output

# Generation with autoencoders?

- Can we use an autoencoder to generate data given a code?
- Yes, but...how to sample the code *z* first?

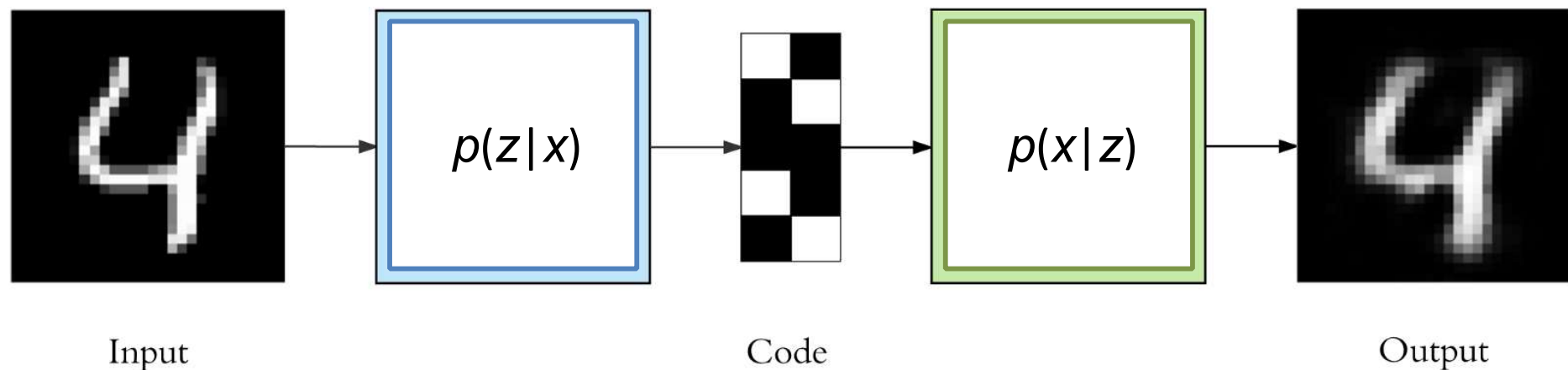*p(z)* ?

Encoder → | Code | → $p(x|z)$ → 

Input | Code | Output

# Generation with autoencoders?

- Can we use an autoencoder to generate data given a code?
- Yes, but...how to sample the code *z* first?



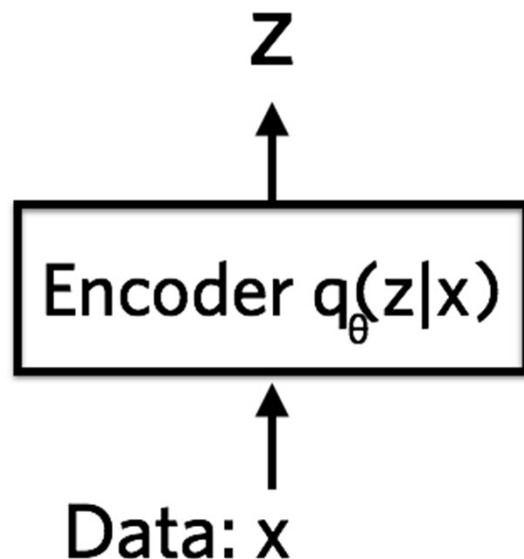Input      Code      Output

In the diagram: $p(z|x)$, $p(x|z)$

# Variational autoencoder

- Let's have a closer look at our (variational) autoencoder components

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



The Encoder is a neural network that takes a data point in input and outputs a hidden representation **z**, usually lower dimensional

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



The Encoder is a neural network that takes a data point in input and outputs a hidden representation **z**, usually lower dimensional

More precisely, the Encoder outputs the parameters of a Gaussian probability density $q_\theta(z|x)$
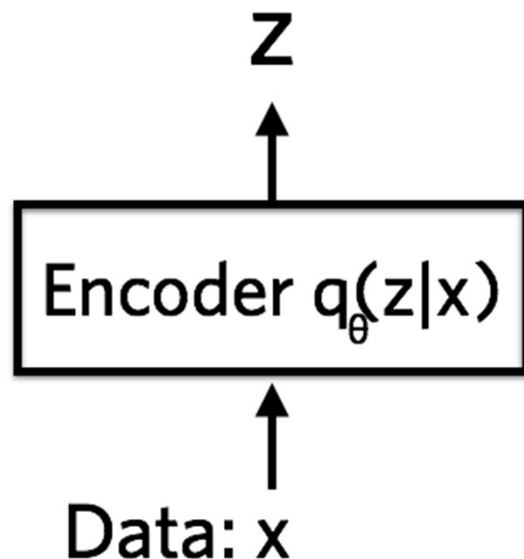
# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components
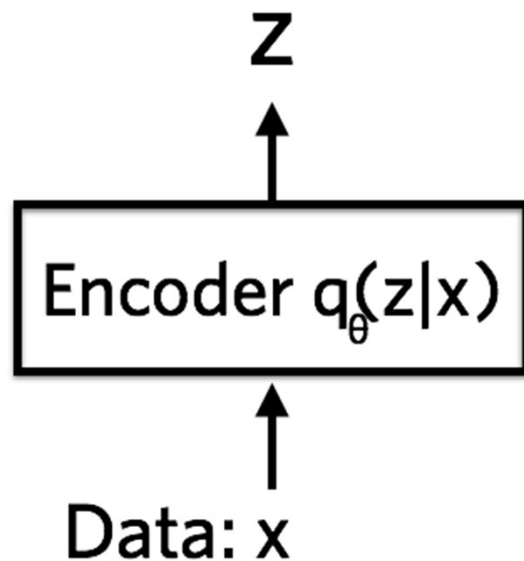


The Encoder is a neural network that takes a data point in input and outputs a hidden representation **z**, usually lower dimensional

More precisely, the Encoder outputs the parameters of a Gaussian probability density $q_\theta(z|x)$

Then, we can sample a representation/code **z** from $q_\theta(z|x)$

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components

The Decoder is another neural network that takes the representation **z** in input and outputs the parameters of $p_\phi(x|z)$, a probability density from which we can sample the data

z

↓

Decoder $p_\phi(x|z)$

↓

Reconstruction: x̃

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components

Z

Z

$$\text{Encoder } q_\theta(z|x)$$

$$\text{Decoder } p_\phi(x|z)$$

Data: x

Reconstruction: x̃

Input: 28 x 28 image of a digit
(784-dimensional binary image)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



$$\text{Encoder } q_\theta(z|x)$$

$$\text{Decoder } p_\phi(x|z)$$

Data: x

Reconstruction: x̃

0 0 0 1
0 0 1 1
0 0 1 0 …

Input: 28 x 28 image of a digit
(784-dimensional binary image)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



z

Encoder $q_\theta(z|x)$

Data: x

z

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

The input is fed to a network (the encoder) with its weights and biases ($\theta$)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



$$Z \qquad\qquad Z$$

Encoder $q_\theta(z|x)$      Decoder $p_\phi(x|z)$

Data: $x$      Reconstruction: $\tilde{x}$

The network outputs mean
and variance for $q_\theta(z|x)$

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Encoder $q_\theta(z|x)$

Decoder $p_\phi(x|z)$

Data: x

Reconstruction: $\tilde{x}$

Sample **z** from **$q_\theta(z|x)$**, where the subscript **θ** means that the probability is parametrised by **θ**

# VAEs: encoder

# VAEs: encoder

Output
$\mu$
$[0.1, 1.2, 0.2, 0.8,...]$

Output
$\sigma$
$[0.2, 0.5, 0.8, 1.3,...]$

Intermediate
$X$
$[X_1 \sim N(0.1, 0.2^2), X_2 \sim N(1.2, 0.5^2), X_3 \sim N(0.2, 0.8^2), X_4 \sim N(0.8, 1.3^2),....]$

↓ sample

Sampled
vector
$[0.28, 1.65, 0.92, 1.98,...]$

# VAEs: neural networks perspective

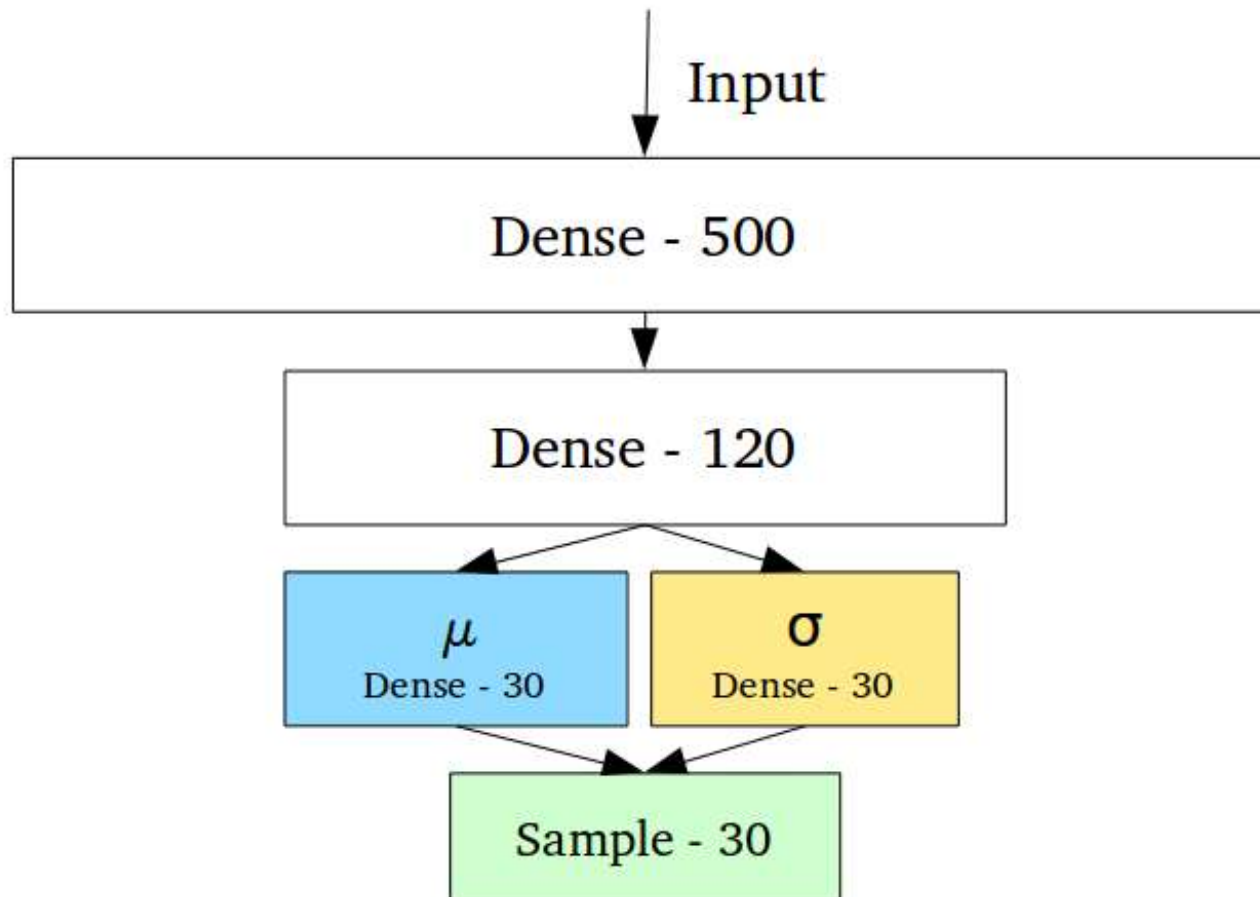- Let's have a closer look at our (variational) autoencoder components

0.1 -0.2 2.2
1.1  0.1  -1 …

**z**

**z**

Encoder $q_\theta(z|x)$

Decoder $p_\phi(x|z)$

Data: x

Reconstruction: x̃

Note that different samples from $q_\theta(z|x)$ will in general yield different codes **z**

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components

-1.1 0.2 0.3
0.1  0.2  2 ...

**z**

Encoder $q_\theta(z|x)$

Data: x

**z**

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

Note that different samples from $q_\theta(z|x)$ will in general yield different codes **z**

# VAEs: neural networks perspective

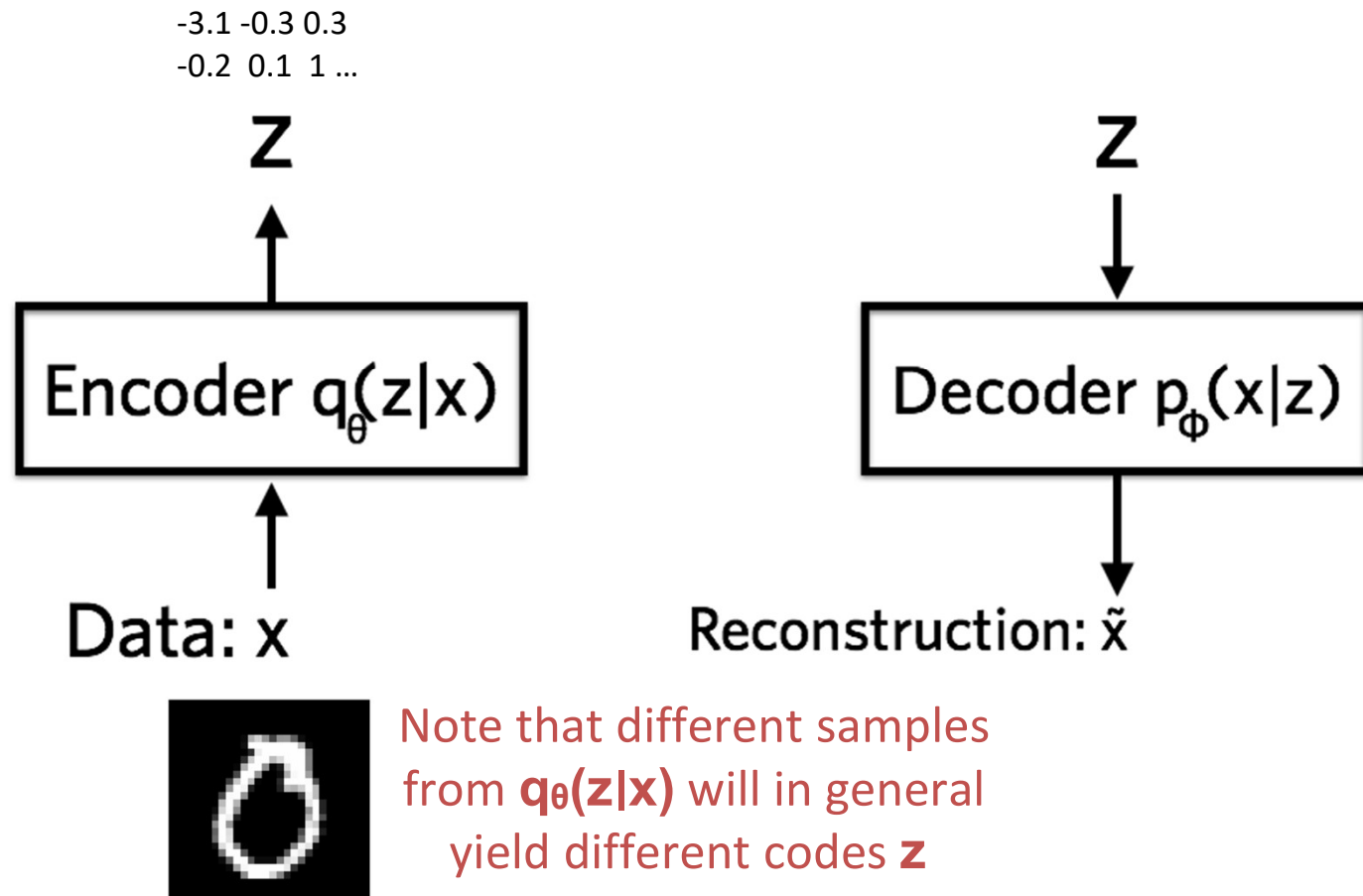- Let's have a closer look at our (variational) autoencoder components

-3.1 -0.3 0.3
-0.2  0.1  1 …

**Z**

**Z**

Encoder $q_\theta(z|x)$

Decoder $p_\phi(x|z)$

Data: x

Reconstruction: x̃

Note that different samples from $q_\theta(z|x)$ will in general yield different codes **z**

# Latent space: AE vs VAEs



Standard Autoencoder
(direct encoding coordinates)
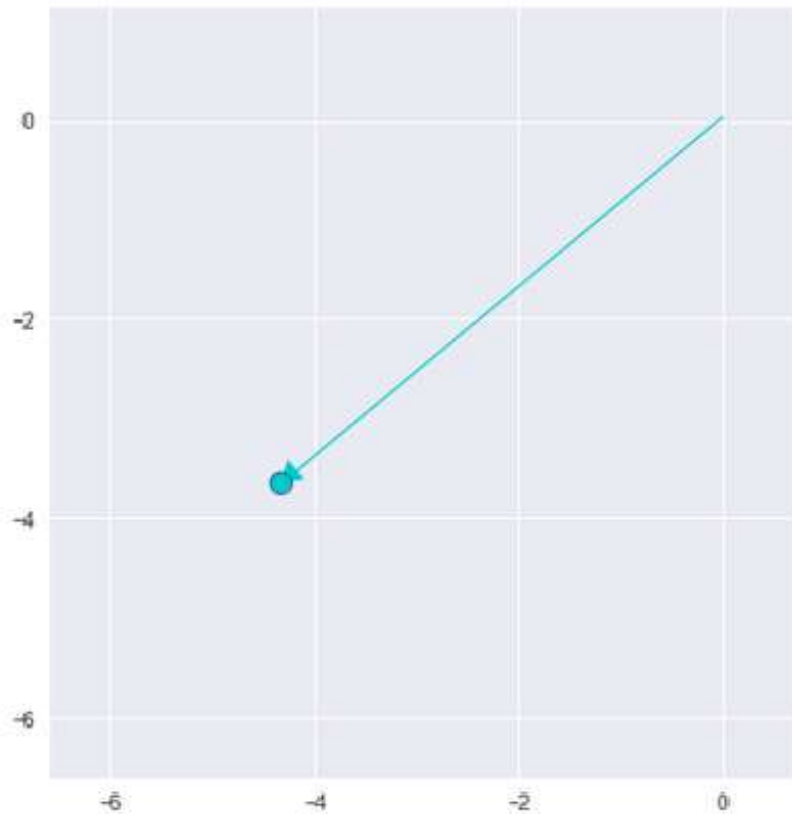
Variational Autoencoder
($\mu$ and $\sigma$ initialize a probability distribution)

# VAEs: neural networks perspective

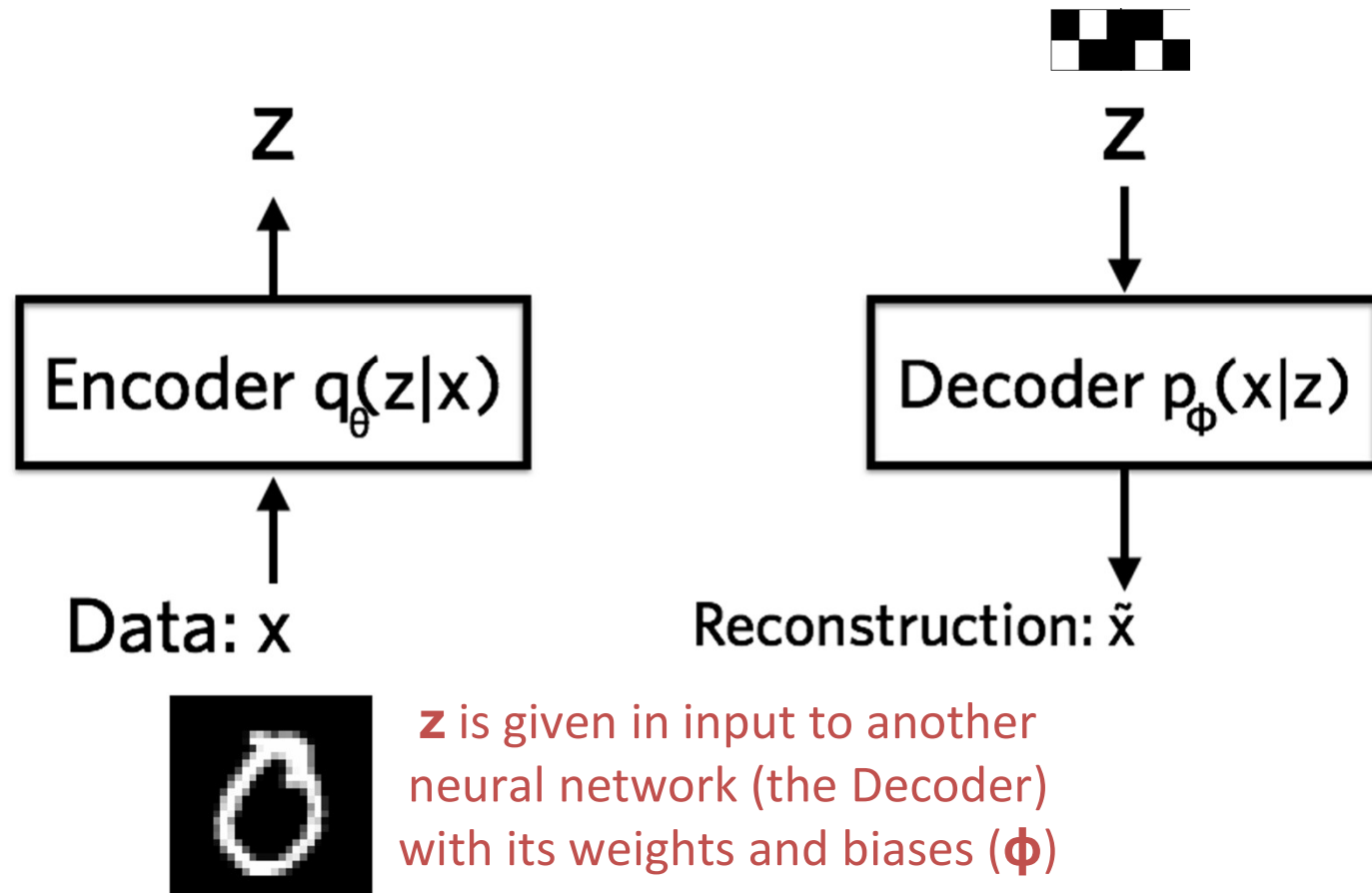- Let's have a closer look at our (variational) autoencoder components



Encoder $q_\theta(z|x)$

Decoder $p_\phi(x|z)$

Data: x

Reconstruction: x̃

**z** is given in input to another neural network (the Decoder) with its weights and biases (**φ**)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Encoder $q_\theta(z|x)$

Data: $x$

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

The network outputs the parameters of the distribution $p_\phi(x|z)$

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Z

Encoder $q_\theta(z|x)$

Data: x

Z

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

Recall that we're dealing with binary images

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Encoder $q_\theta(z|x)$

Data: $x$

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

So we can teach our decoder to output 784 Bernoulli parameters (1 for each pixel)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



$$Z$$

$$\text{Encoder } q_\theta(z|x)$$

Data: x

$$Z$$

$$\text{Decoder } p_\phi(x|z)$$

0 0.1 0.9
0 0.3 0.8
0 0.1 0.8 ...

Reconstruction: x̃

So we can teach our decoder to output 784 Bernoulli parameters (1 for each pixel)

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Z

Encoder $q_\theta(z|x)$

Data: x

Z

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

Finally we sample a new image from the Bernoulli distributions
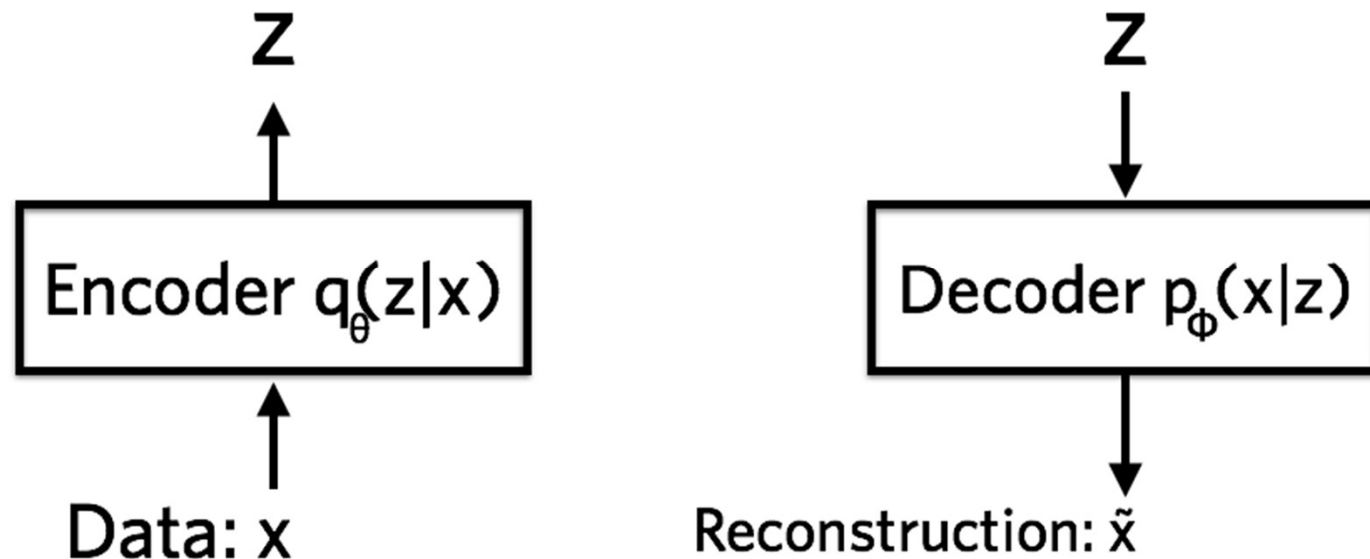
# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



How much information is lost when going from the low-dimensional representation **z** to the higher dimensional reconstructed **x**?
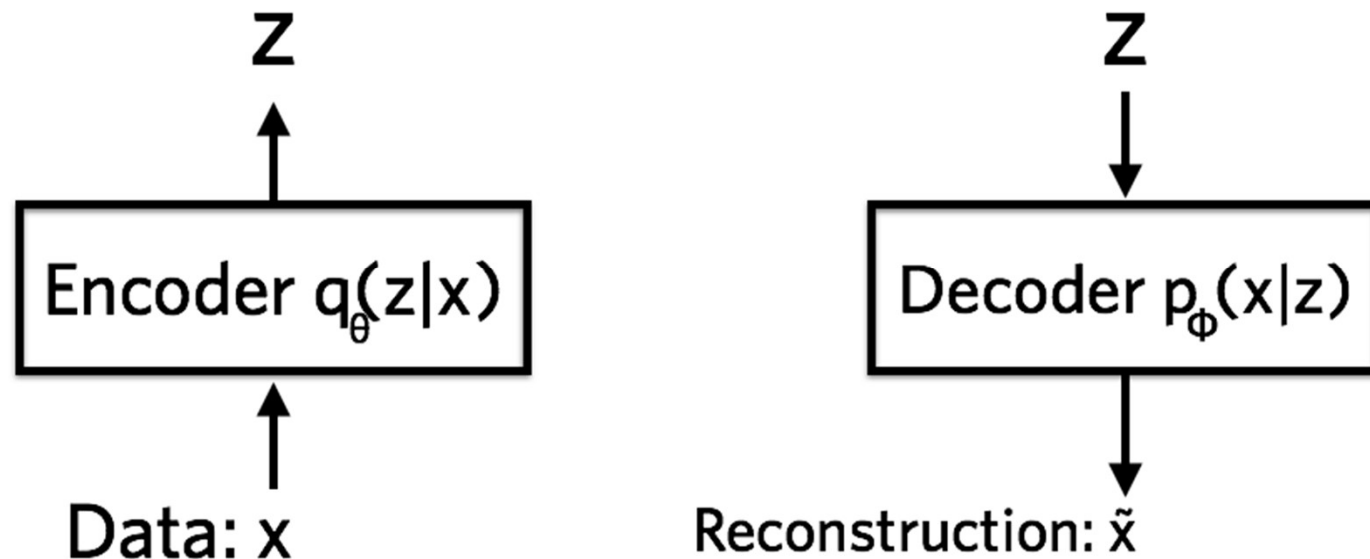
# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



This can be measured using the reconstruction log-likelihood **log $p_\phi$(x|z)**
It tells us how effectively the decoder has learned to reconstruct **x** given **z**
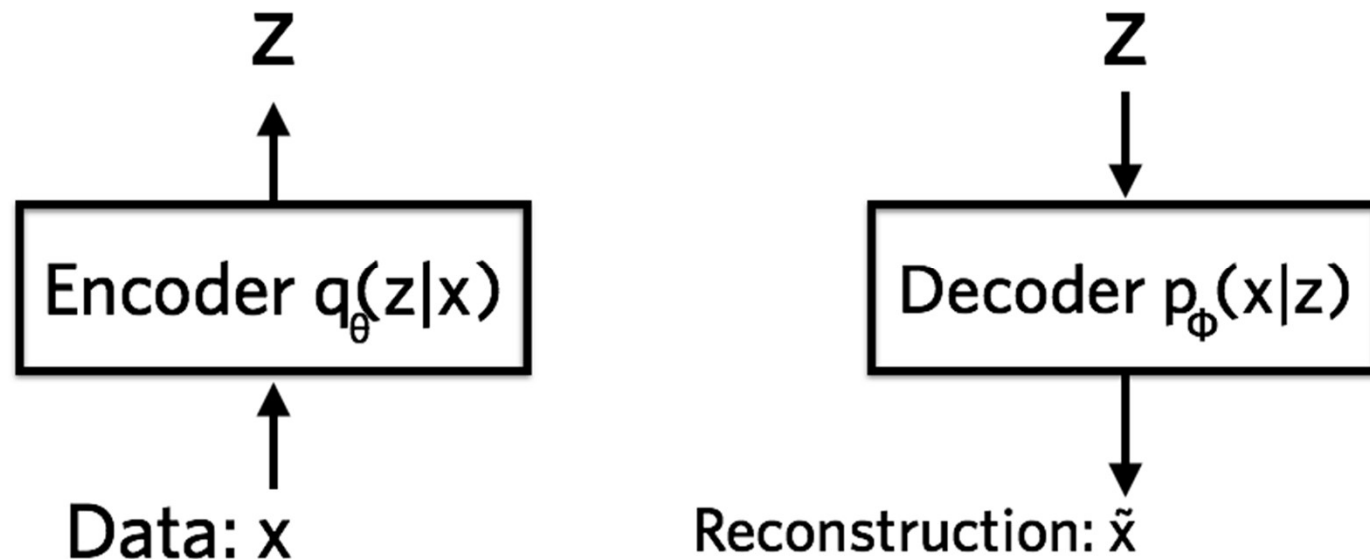
# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



Encoder $q_\theta(z|x)$

Data: $x$

Decoder $p_\phi(x|z)$

Reconstruction: $\tilde{x}$

Sounds nice, but we need a loss if we want to do back-propagation
and learn the optimal parameters of the two networks (encoder & decoder)

# VAEs: neural networks perspective

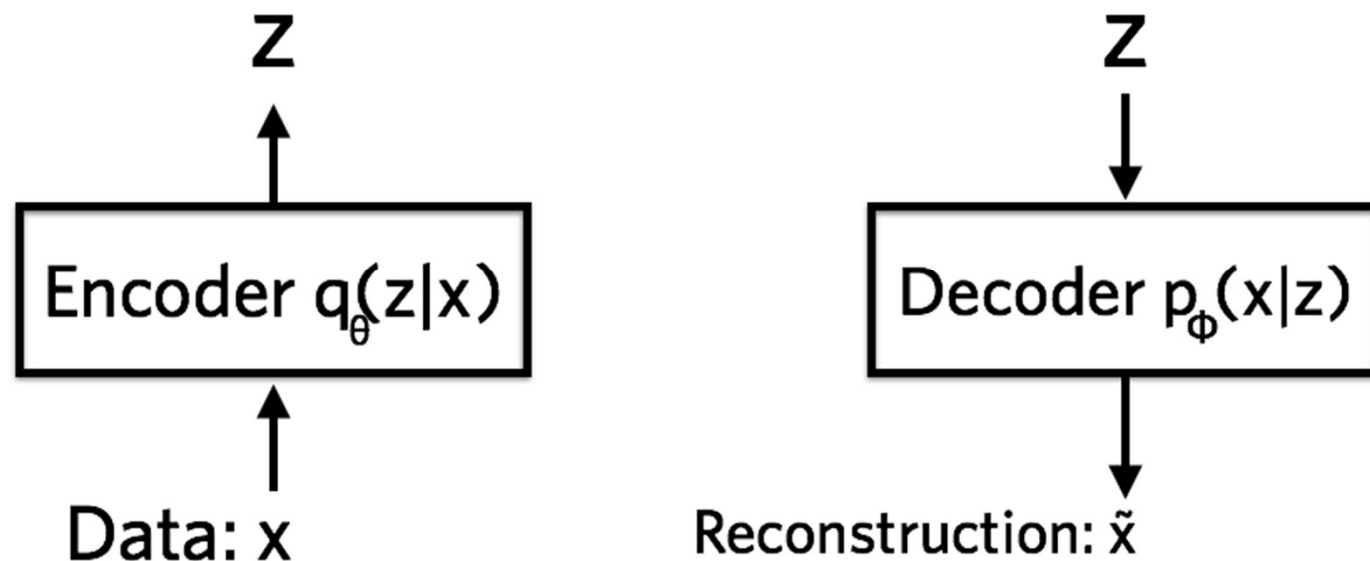- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \parallel p(z))$$

loss function for i-th datapoint - the total loss is the sum of all the $l_i$

# VAEs: neural networks perspective

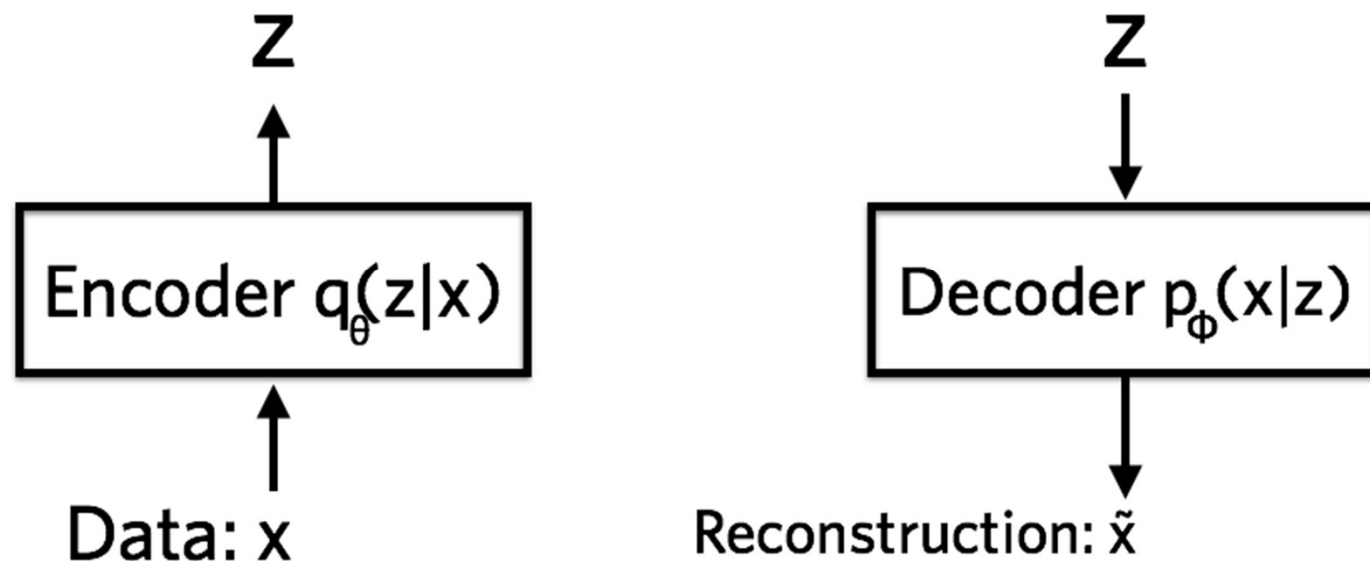- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

Remember, for one data point we can sample many codes **z**

37

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components
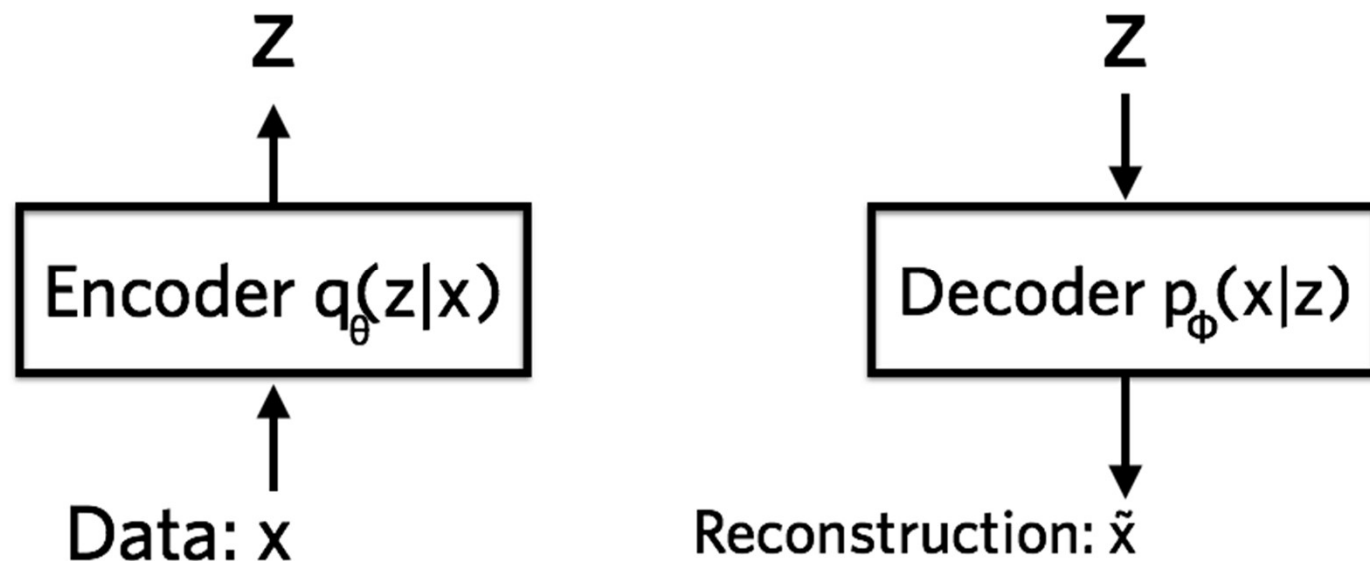


$$l_i(\theta, \phi) = \boxed{-\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)]} + \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

This is why here we have the expected log-likelihood (negative, we're minimising)

# VAEs: neural networks perspective

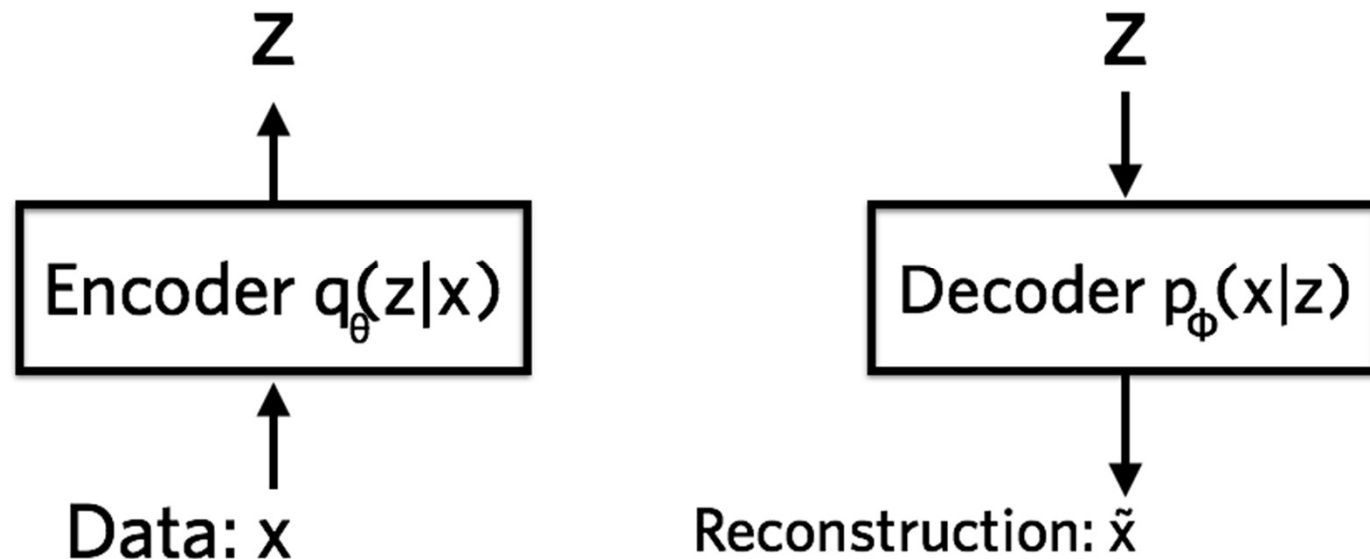- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = \boxed{-\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)]} + \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

This term encourages decoder to learn to reconstruct the data

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components
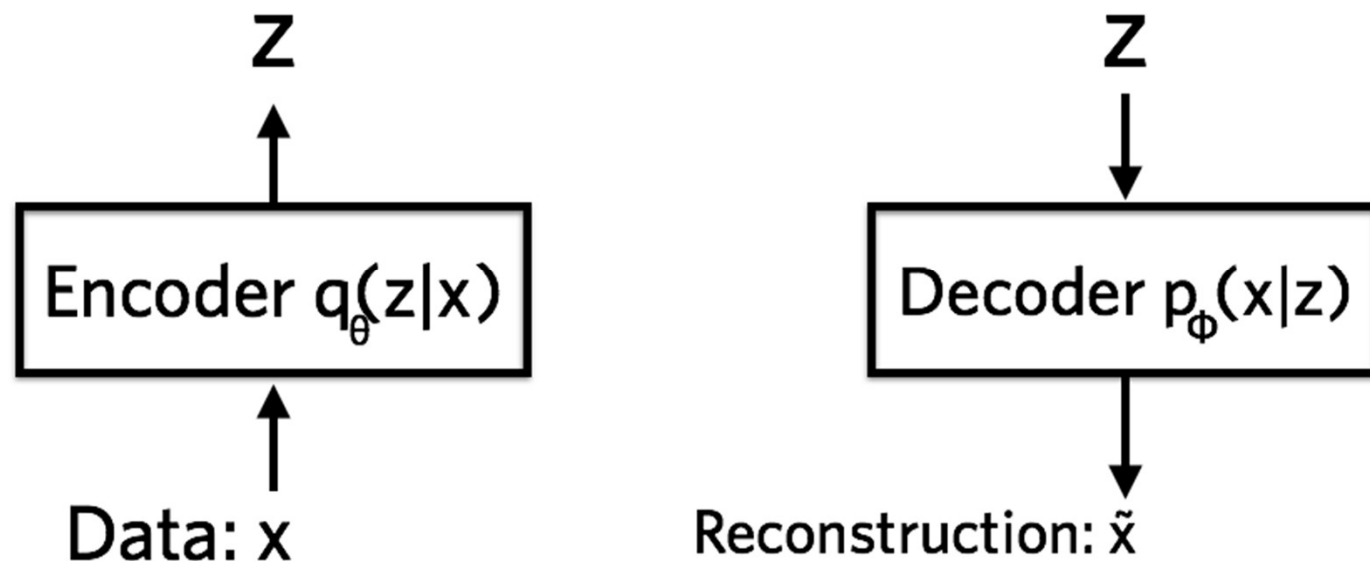


$$l_i(\theta, \phi) = \boxed{-\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)]} + \mathbb{KL}(q_\theta(z \mid x_i) \,\|\, p(z))$$

Bad encoder (e.g., high prob of black pixel when it should be white) = large cost

# VAEs: neural networks perspective

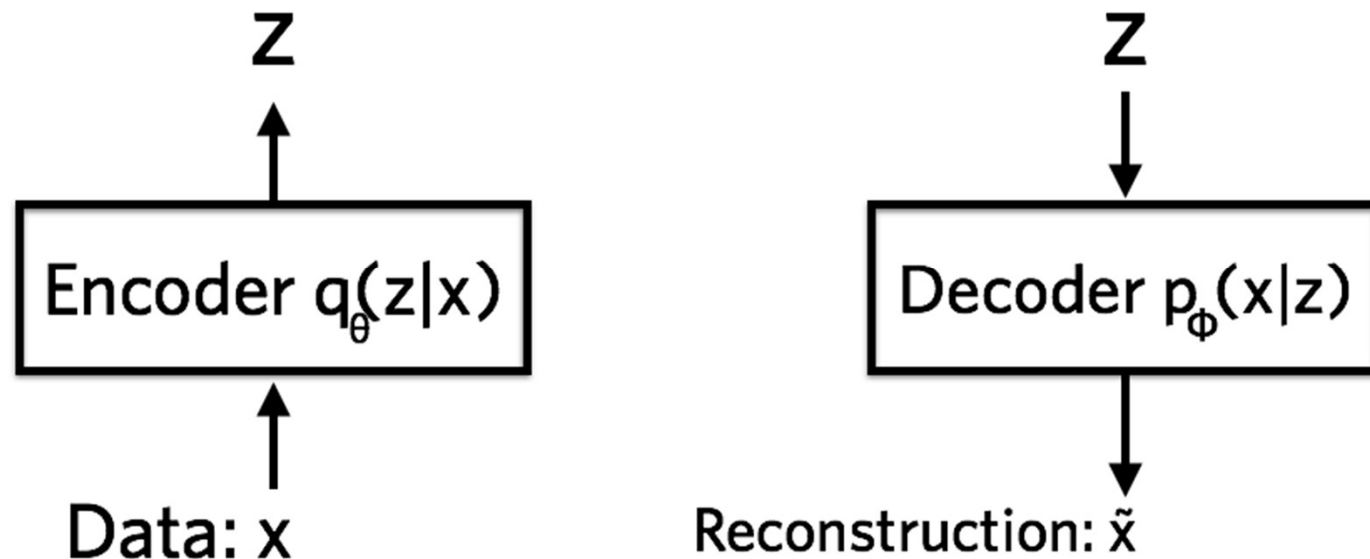- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

Regularisation term - measures how close **q** is to **p**

# VAEs: neural networks perspective

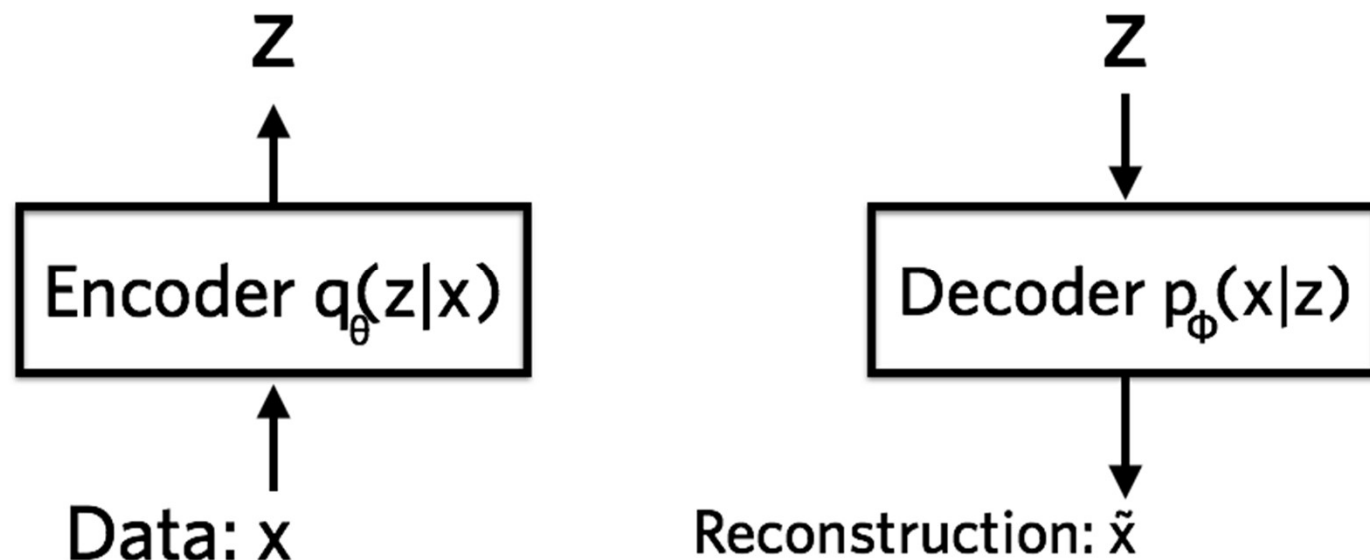- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \parallel p(z))$$

In VAs we let **p(z) = Gaussian(0,1)**

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



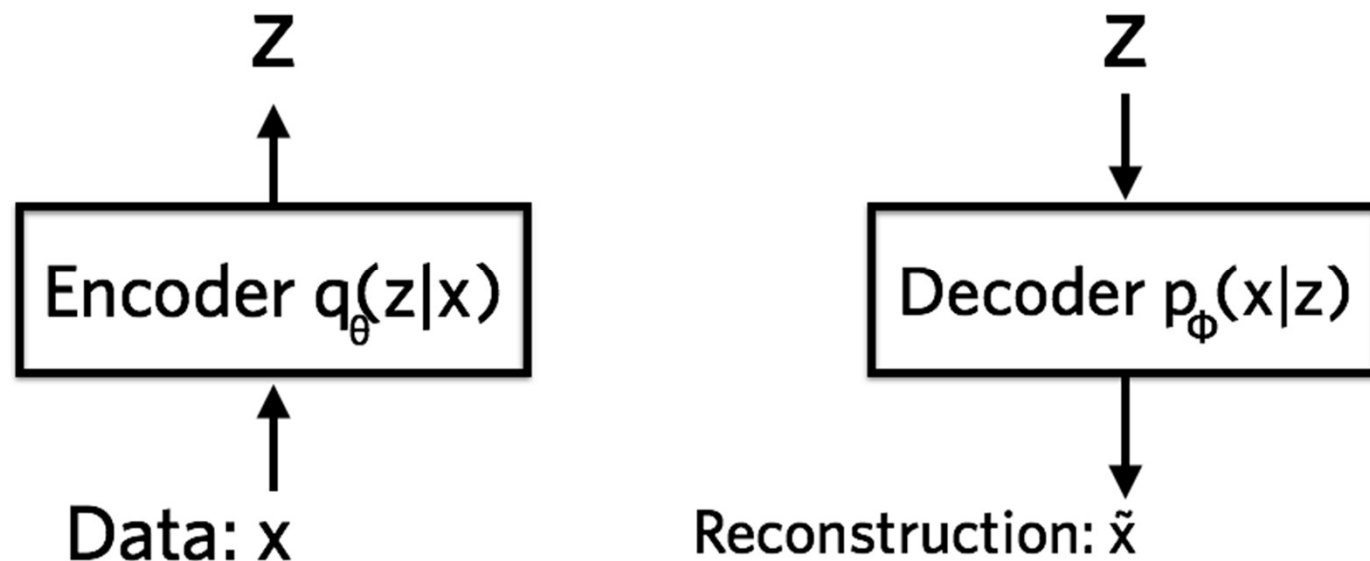$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \parallel p(z))$$

This is added to make sure the encoder doesn't cheat and map each datapoint in different regions of the space

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \;||\; p(z))$$

This is bad because we want e.g. different images of same number to lie close in the embedding space

# VAEs: neural networks perspective

- Let's have a closer look at our (variational) autoencoder components



$$l_i(\theta, \phi) = -\mathbb{E}_{z \sim q_\theta(z|x_i)}[\log p_\phi(x_i \mid z)] + \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

The space has to be "meaningful", so we penalise this behaviour

# Meaning of the regularisation



What we require

What we may inadvertently end up with

# VAEs: neural networks perspective

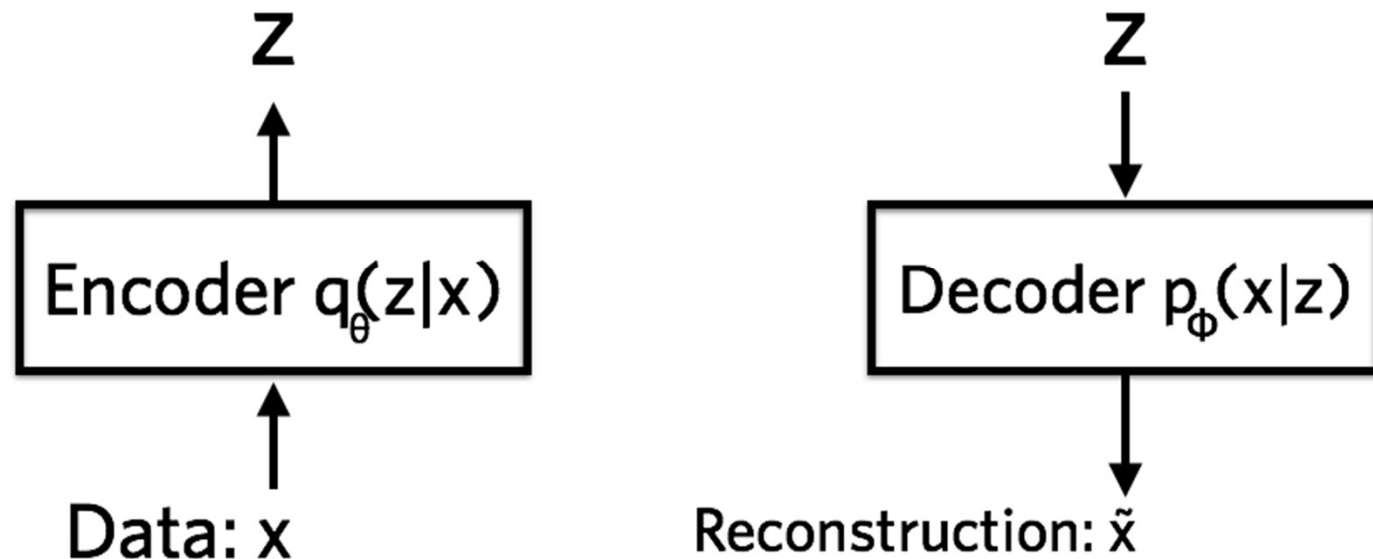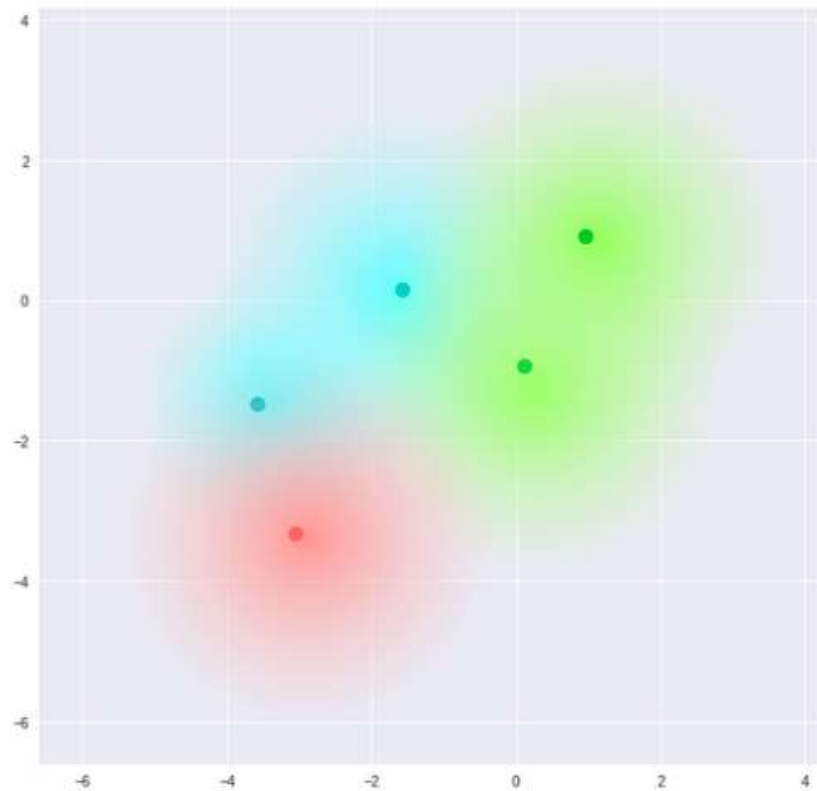- Let's have a closer look at our (variational) autoencoder components



Given this architecture, we can use back-propagation to compute the gradients of the loss wrt the networks parameters and then optimise using any variant of gradient descent

# VAEs architecture



Prior distribution: $p_\theta(\mathbf{z})$

$\mathbf{z}$-space

Encoder: $q_\varphi(\mathbf{z}|\mathbf{x})$

Decoder: $p_\theta(\mathbf{x}|\mathbf{z})$

$\mathbf{x}$-space

Dataset: $\mathbf{D}$

# VAEs architecture

# VAEs: latent space exploration



**Want**: find images between "face with glasses" to "face without glasses"

First: Find representation of faces in latent space (by giving it as input of encoder)

# VAEs: latent space exploration

Glasses



**What we do**: In the latent space, you computed the translation vector from "face without glasses" to "face with glasses" from the embeddings of two faces, without and with glasses

# VAEs: latent space exploration

Glasses



Add the translation vector to the latent representation then
decode this representation to map it back to the image space



Face with glasses

Glasses

Face without glasses

# VAEs: latent space exploration



Latent space interpolation

# VAEs: neural networks perspective

- ...but why exactly is this called a "variational" autoencoder?

# VAEs: probabilistic perspective

- To understand why, we need to look at variational autoencoders from a different perspective, i.e., that of a probability model

# VAEs: probabilistic perspective

- To understand why, we need to look at variational autoencoders from a different perspective, i.e., that of a probability model

- Let there be a generative model of the data **x** and the latent variables **z** with joint probability

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

# VAEs: probabilistic perspective

- To understand why, we need to look at variational autoencoders from a different perspective, i.e., that of a probability model

- Let there be a generative model of the data **x** and the latent variables **z** with joint probability
**p(x,z) = p(x|z)p(z)**

- First we sample **z** from prior **p(z)**

- Then we sample **x** from the likelihood **p(x|z)**

# VAEs: probabilistic perspective

- In this context, learning is called *inference*
- We want to infer the optimal parameters of **p(x)**
- In other words, we want to maximise **p(x)**
- **log p(x) = log p(x,z)/p(z|x)**, to be precise

# VAEs: probabilistic perspective

- In this context, learning is called *inference*

- We want to infer the optimal parameters of $p(x)$

- In other words, we want to maximise $p(x)$

- $\log p(x) = \log p(x,z)/p(z|x)$, to be precise

- But $p(z|x) = p(x,z)/p(x)$, and computing $p(x)$ takes an exponential time since $p(x) = \int p(x|z)p(z)dz$

- $p(z|x)$ is called the posterior and its often intractable in this type of problems...

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior $p(z|x)$ with a family of distributions $q_\lambda(z|x)$

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\mathbb{KL}(q_\lambda(z \mid x) \mid\mid p(z \mid x)) =$$

$$\mathbf{E}_q[\log q_\lambda(z \mid x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

$$\mathbb{KL}(q_\lambda(z \mid x) \| p(z \mid x)) =$$

$$\mathbf{E}_q[\log q_\lambda(z \mid x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

$$KL\left(q_\lambda(z \mid x) \| P(z \mid x)\right) = E_q\left(\log \frac{q_\lambda(z \mid x)}{P(z \mid x)}\right)$$

$$= E_q\left(\log q_\lambda(z \mid x) - \log P(z \mid x)\right)$$

$$= E_q\left(\log q_\lambda(z \mid x)\right) - E_q\left(\log \frac{P(z, x)}{P(x)}\right) \quad \text{Constant}$$

$$= E_q\left(\log q_\lambda(z \mid x)\right) - E_q\left(\log P(x, z)\right) + E_q\left(\log P(x)\right)$$

$$= E_q\left(\log q_\lambda(z \mid x)\right) - E_q\left(\log P(x, z)\right) + \log P(x)$$

$$KL \geq 0 \implies \log P(x) \geq \underline{E_q\left(\log P(x, z)\right) - E_q\left(\log q_\lambda(z \mid x)\right)}$$

ELBO (Evidence lower bound)

if want to max $\log P(x)$, then max ELBO.

(by Jiangao Zhang.)

$$\text{ELBO} = E_q\left(\log p(x, z)\right) - E_q\left(\log q_\lambda(z|x)\right)$$

$$= E_q\left(\log p(x|z) \cdot p(z)\right) - E_q\left(\log q_\lambda(z|x)\right)$$

$$= E_q\left(\log p(x|z)\right) + E_q\left(\log p(z)\right) - E_q\left(\log q_\lambda(z|x)\right)$$

$$= E_q\left(\log p(x|z)\right) - E_q\left(\log \frac{q_\lambda(z|x)}{p(z)}\right)$$

$$= E_q\left(\log p(x|z)\right) - KL\left(q_\lambda(z|x) \,\|\, p(z)\right)$$

---

*     objective function of VAE.

*     $p(z)$ is the prior, normally the Gaussian. Distribution

x   The first term is the maximum likelihood of the Decoder.

(by Jiangeo Zhang)

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\mathbb{KL}(q_\lambda(z \mid x) \mid\mid p(z \mid x)) =$$

$$\mathbf{E}_q[\log q_\lambda(z \mid x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

We want this to be small

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$q_\lambda^*(z \mid x) = \arg\min_\lambda \mathbb{KL}(q_\lambda(z \mid x) \mid\mid p(z \mid x))$$

This is what we're looking for

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$q_\lambda^*(z \mid x) = \arg\min_\lambda \mathbb{KL}(q_\lambda(z \mid x) \;||\; p(z \mid x))$$

But this isn't good…

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\mathbb{KL}(q_\lambda(z \mid x) \,||\, p(z \mid x)) =$$

$$\mathbf{E}_q[\log q_\lambda(z \mid x)] - \mathbf{E}_q[\log p(x, z)] + \log p(x)$$

Let us introduce this function

$$ELBO(\lambda) = \mathbf{E}_q[\log p(x, z)] - \mathbf{E}_q[\log q_\lambda(z \mid x)]$$

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\log p(x) = ELBO(\lambda) + \mathbb{KL}(q_\lambda(z \mid x) \, || \, p(z \mid x))$$

Then we can write this

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\log p(x) = ELBO(\lambda) + \mathbb{KL}(q_\lambda(z \mid x) \mid\mid p(z \mid x))$$

KL is always >= 0, so to minimise KL we can maximise ELBO!

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$\log p(x) = ELBO(\lambda) + \mathbb{KL}(q_\lambda(z \mid x) \mid\mid p(z \mid x))$$

Maximising ELBO means 1) **q** close to **p** and 2) higher **p** (better generator)

$$ELBO_i(\theta, \phi) = \mathbb{E}q_\theta(z \mid x_i)[\log p_\phi(x_i \mid z)] - \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

# VAEs: probabilistic perspective

- Which is where the **variational** elements comes in!

- Variational inference approximates the posterior **p(z|x)** with a family of distributions **q$_\lambda$(z|x)**

- Of course we want our approximation to be good, i.e., close to the true posterior

$$ELBO(\lambda) = \mathbf{E}_q[\log p(x, z)] - \mathbf{E}_q[\log q_\lambda(z \mid x)]$$

$$ELBO_i(\theta, \phi) = \mathbb{E}q_\theta(z \mid x_i)[\log p_\phi(x_i \mid z)] - \mathbb{KL}(q_\theta(z \mid x_i) \mid\mid p(z))$$

We link the network and probabilistic perspective by explicating the parameters of q and p and noting that the above is the (negative of the) loss function of a variational autoencoder

# VAEs on the Web

- Various visualisations: https://jaan.io/what-is-variational-autoencoder-vae-tutorial/

- Interactive VAEs: https://www.siarez.com/projects/variational-autoencoder

- Morphing faces: http://vdumoulin.github.io/morphing_faces/online_demo.html

- MNIST demo: http://dpkingma.com/sgvb_mnist_demo/demo.html

# VAEs on the Web

- Various visualisations: https://jaan.io/what-is-variational-autoencoder-vae-tutorial/

- Interactive VAEs: https://www.siarez.com/projects/variational-autoencoder

- Morphing faces: http://vdumoulin.github.io/morphing_faces/online_demo.html

- MNIST demo: http://dpkingma.com/sgvb_mnist_demo/demo.html

- Chemical molecule design with VAEs: https://github.com/aspuru-guzik-group/chemical_vae

- Next Topic
  - Generative adversarial networks