

# Deep Generative Models

Ulrich Paquet  
DeepMind

Transylvanian Machine Learning Summer School  
16-22 July 2018, Cluj-Napoca, Romania

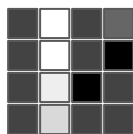
See also

<https://www.youtube.com/watch?v=xTsnNcctvmU>

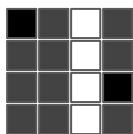
for a recording of a (very, *very*) similar explanation!

## 2-minute exercise

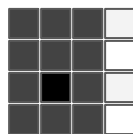
Talk to your friend next to you, and tell him or her everything you can about this data set:



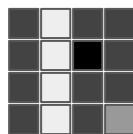
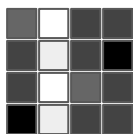
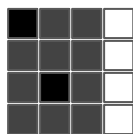
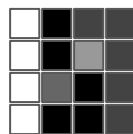
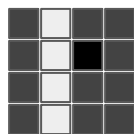
$x^{(1)}$



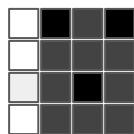
$x^{(2)}$



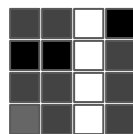
...



...

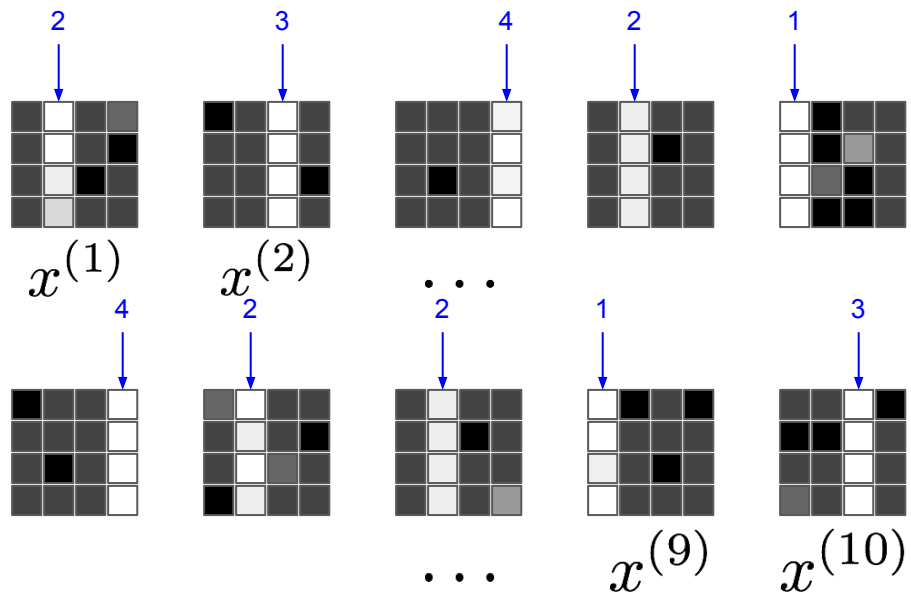


$x^{(9)}$



$x^{(10)}$

# Data



# Data manifold

We can capture most of the variability in the data through **one** number

$$z^{(n)} = 1 \text{ or } 2, 3, 4$$

for each image  $n$ , even though each image is 16 dimensional

How?

# How?

1. Take  $z^{(n)} = 2$
2. Draw bar in column 2 of image
3. Et voila! You have  $x^{(n)}$

$$z^{(n)} = 2$$

Some bar-drawing process

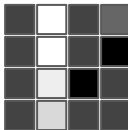
$$x^{(n)}$$

# How?

1. Take  $z^{(n)} = 2$
2. Draw bar in column 2 of image
3. Et voila! You have  $x^{(n)}$

$$z^{(n)} = 2$$

Maybe some neural network, that takes  $z$  as input, and outputs a 16-dimensional vector  $x$ ...?

$$x^{(n)}$$


## 3-minute exercise

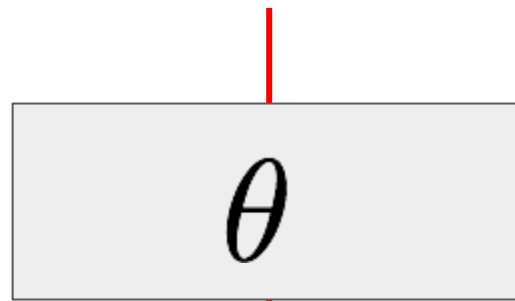
Write or draw a function (like a multi-layer perceptron) that takes  $z \in \mathbb{R}$  and produces  $x$

Is your input one-dimensional?

Is your output 16-dimensional?

Identify all the “tunable” parameters  $\theta$  of your function

$$z^{(n)} = 2$$



$$x^{(n)}$$

A 4x4 grid of squares, representing a 16-dimensional output vector. The grid contains a mix of white, light gray, and dark gray squares.



## 3-minute exercise

Write or draw a function (like a multi-layer perceptron) that takes  $z \in \mathbb{R}$  and produces  $\mathcal{X}$

Is your input one-dimensional?

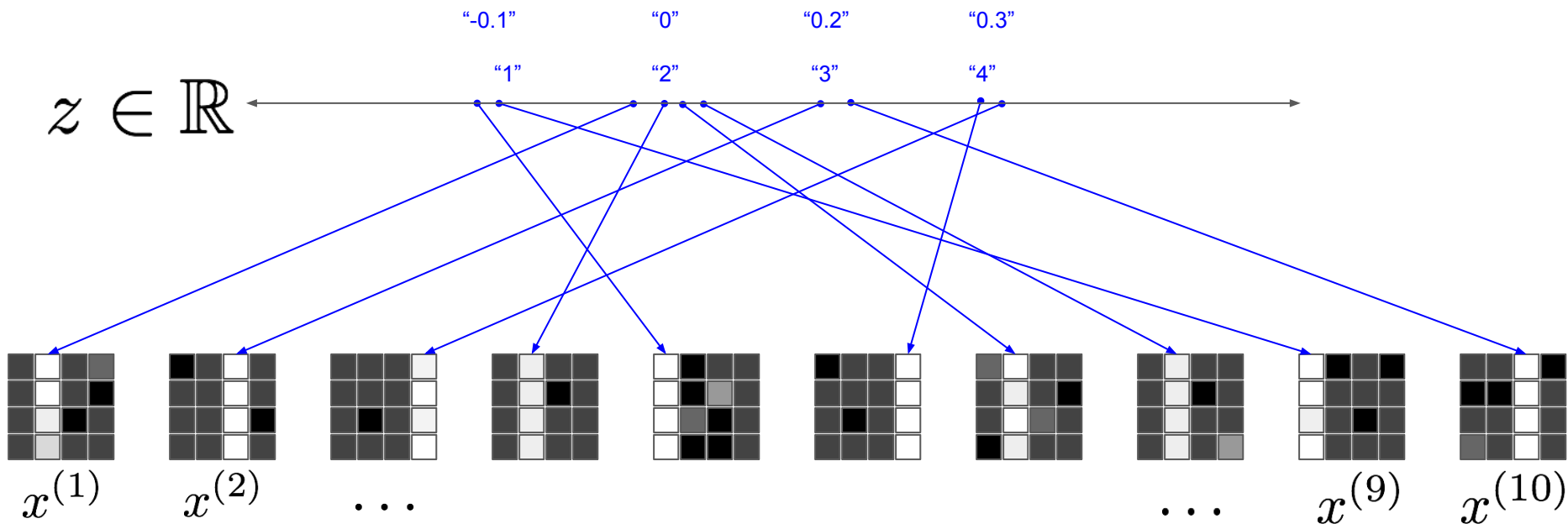
Is your output 16-dimensional?

Identify all the “tunable” parameters  $\theta$  of your function

scratch space

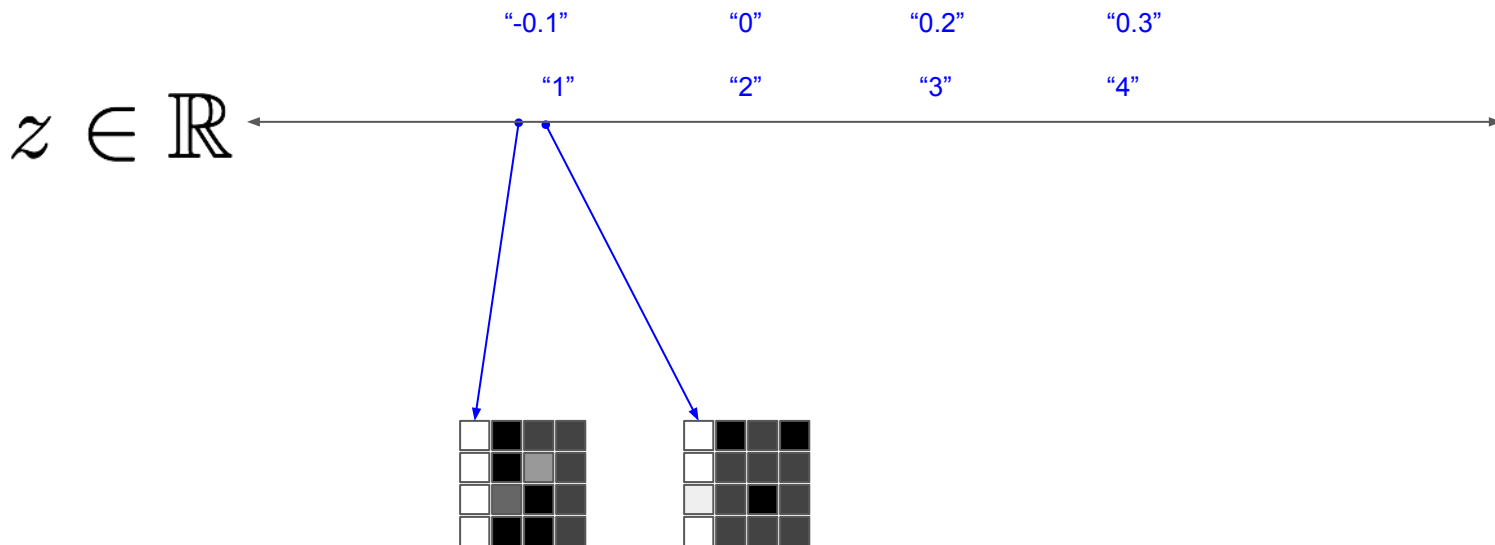
# Data manifold

The 16-dimensional images live on a 1-dimensional manifold, plus some “noise”



## ...and noise

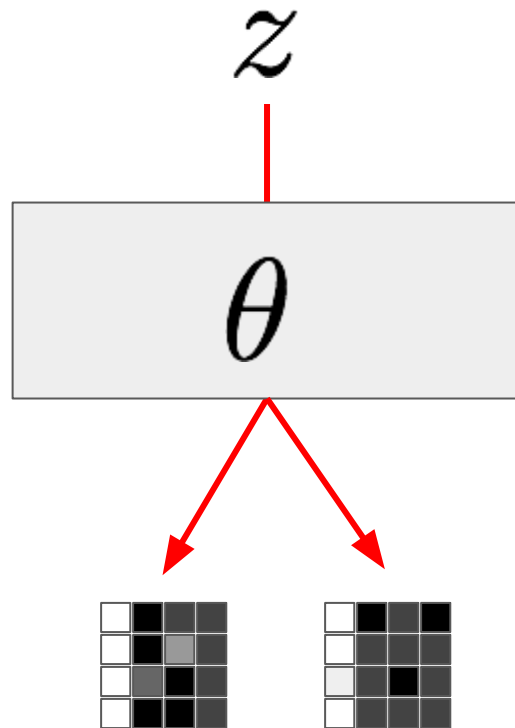
The 16-dimensional images live on a 1-dimensional manifold, plus some “noise”



## 3-minute exercise

Change your multi-layer perceptron to take  $z$  and produce a distribution over  $x$

$$p_{\theta}(x|z)$$



## 3-minute exercise

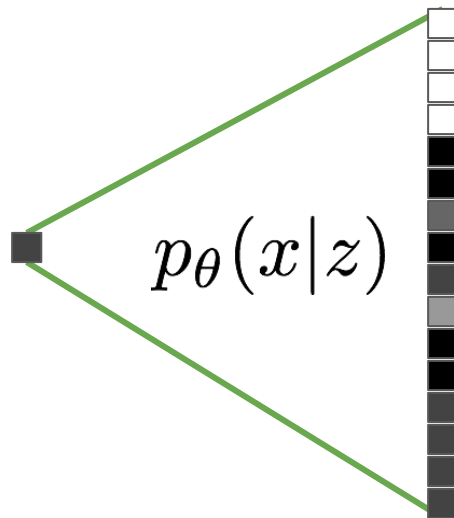
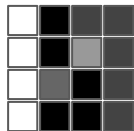
Change your multi-layer perceptron to take  $z$  and produce a distribution over  $x$

$$p_{\theta}(x|z)$$

scratch space

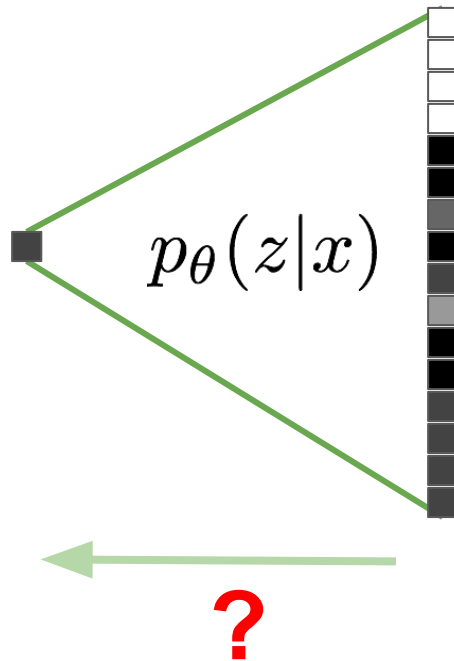
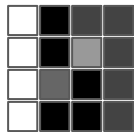
# Decoder

```
def generative_network(z, ...):  
    ...  
    return bernoulli_logits    # for binary pixels
```



$$z \rightarrow x$$

# Inference



# Inversing our world

Two BIG problems to solve:

## Inference

You wrote down  $p_{\theta}(x|z)$  and can compute it.

Say I give you  $x$ . Keeping  $\theta$  fixed, what was  $z$ ? Or  $p_{\theta}(z|x)$ ?

## Learning

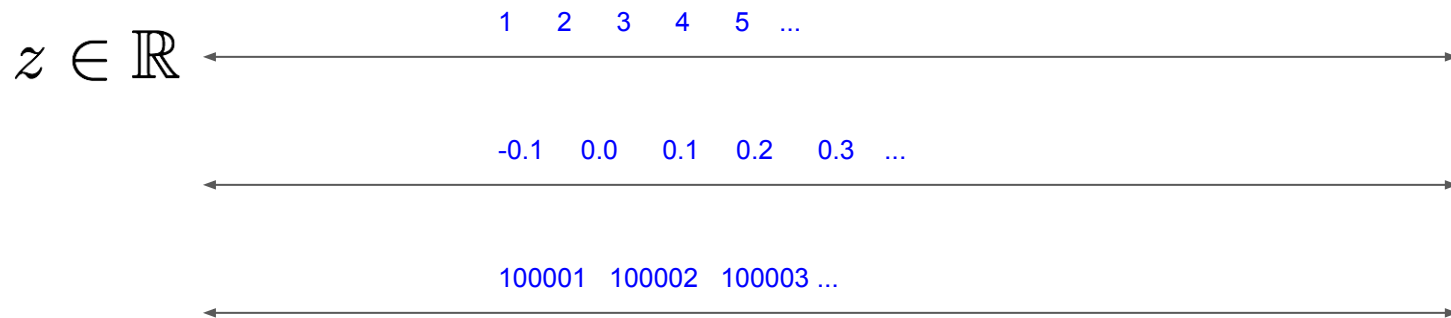
Is there a better (best)  $\theta$  to generate the **observed**  $x$  from  $z$ ?



# Inference

You wrote down  $p_{\theta}(x|z)$  and can compute it.

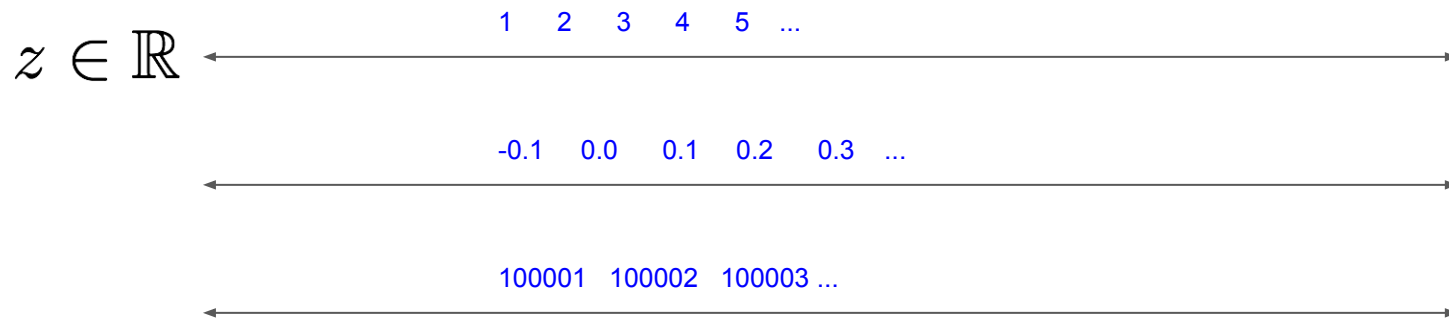
Say I give you  $x$ . Keeping  $\theta$  fixed, what was  $z$ ? Or  $p_{\theta}(z|x)$ ?



# Inference

You wrote down  $p_{\theta}(x|z)$  and can compute it.

Say I give you  $x$ . Keeping  $\theta$  fixed, what was  $z$ ? Or  $p_{\theta}(z|x)$ ?



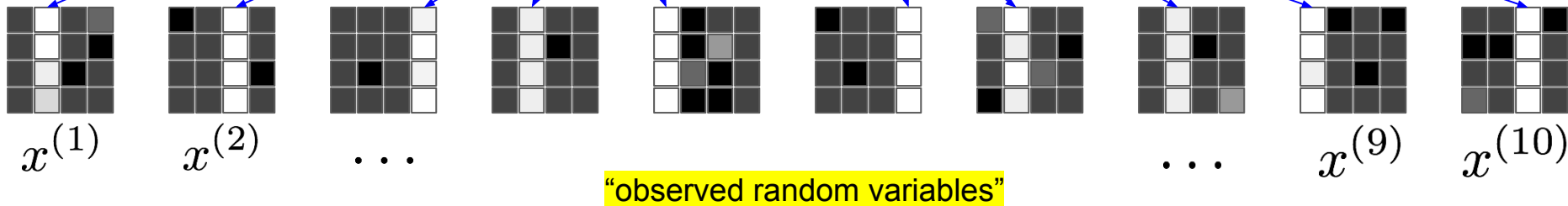
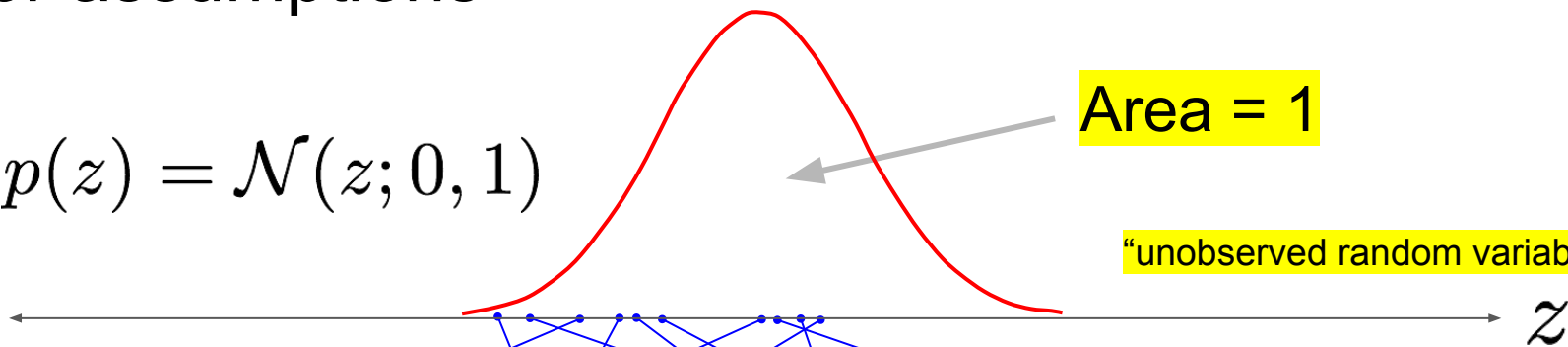
**To really answer that question, we need some notion of where we might have started! No inference without prior assumptions :)**

# Prior assumptions

$$p(z) = \mathcal{N}(z; 0, 1)$$

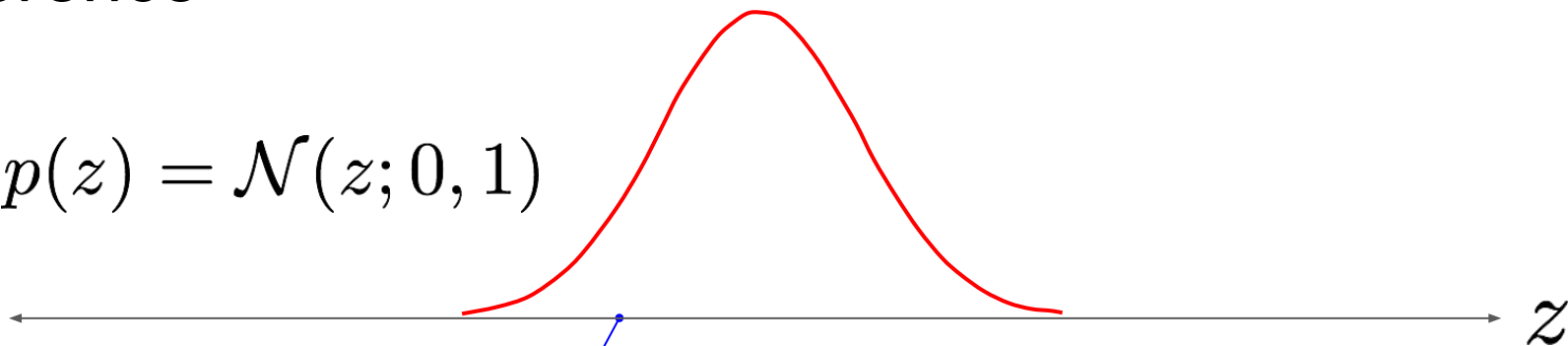
Area = 1

“unobserved random variables”

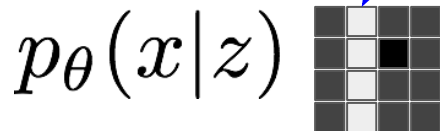


# Inference

$$p(z) = \mathcal{N}(z; 0, 1)$$

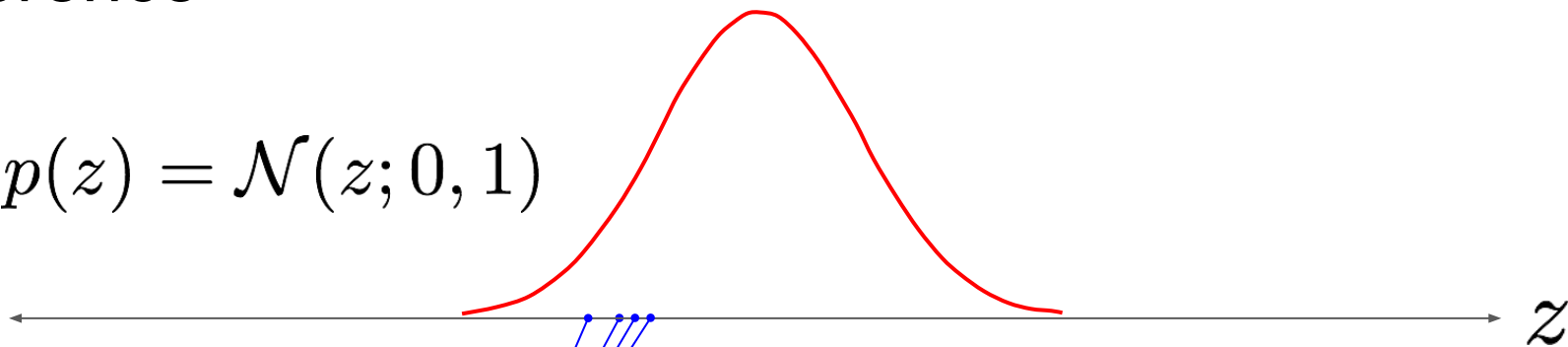


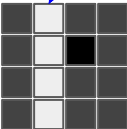
I give you  $x$ . Keeping  $\theta$  fixed, what was  $z$ ?



# Inference

$$p(z) = \mathcal{N}(z; 0, 1)$$

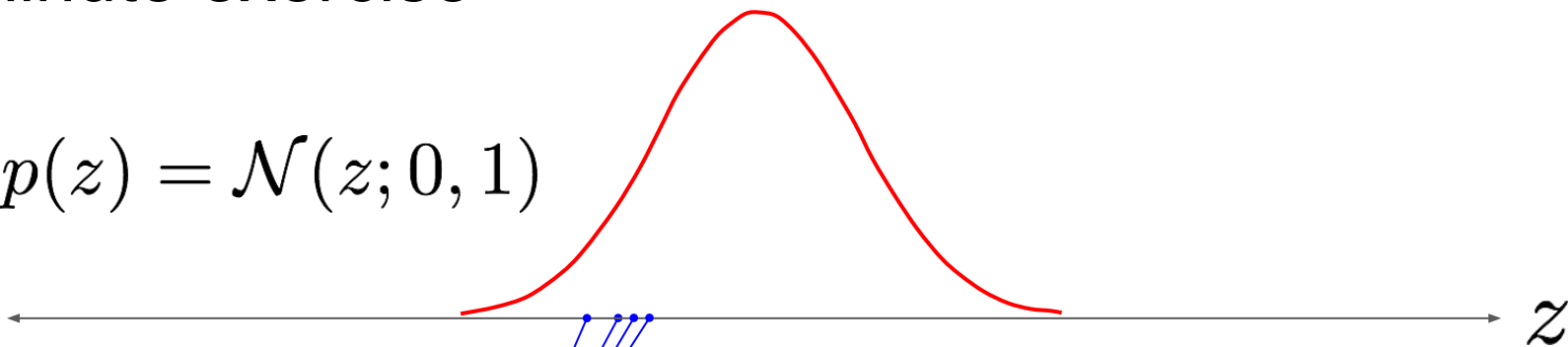


$$p_\theta(x|z)$$


I give you  $x$ . Keeping  $\theta$  fixed, what was  $z$ ?

## 3-minute exercise

$$p(z) = \mathcal{N}(z; 0, 1)$$



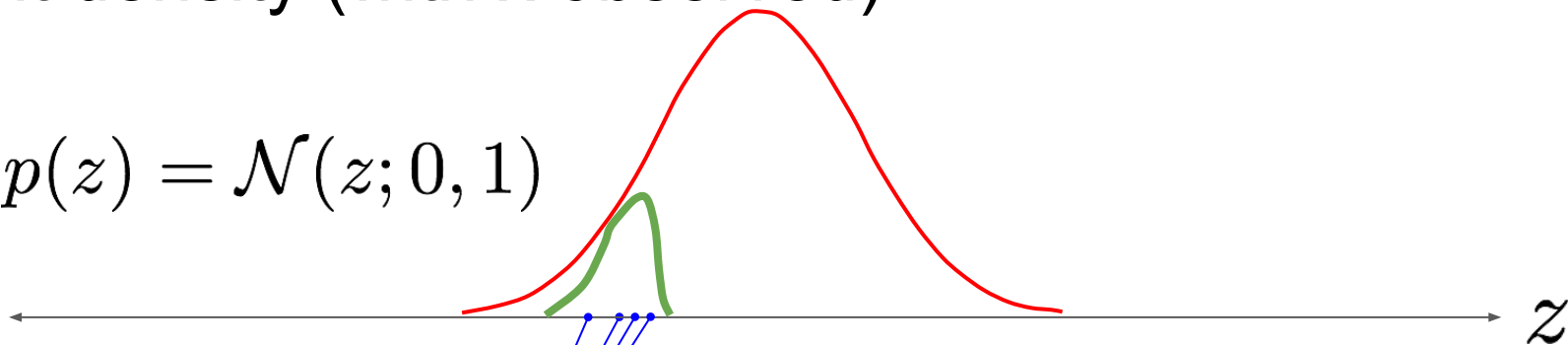
Assuming the largest value of  $p_\theta(x|z)$  is 1,  
draw

$$p_\theta(x, z) = p_\theta(x|z) p(z)$$

as a function of  $z$  on the same axis as above

# Joint density (with x observed)

$$p(z) = \mathcal{N}(z; 0, 1)$$

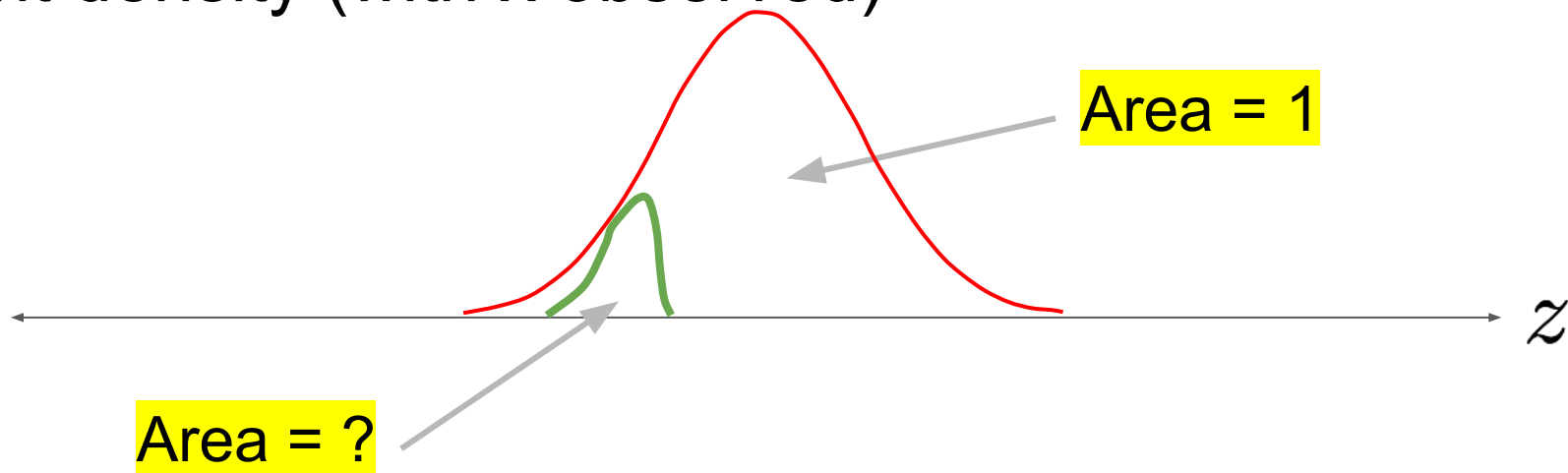


Assuming the largest value of  $p_\theta(x|z)$  is 1,  
draw

$$p_\theta(x, z) = p_\theta(x|z) p(z)$$

as a function of  $z$  on the same axis as above

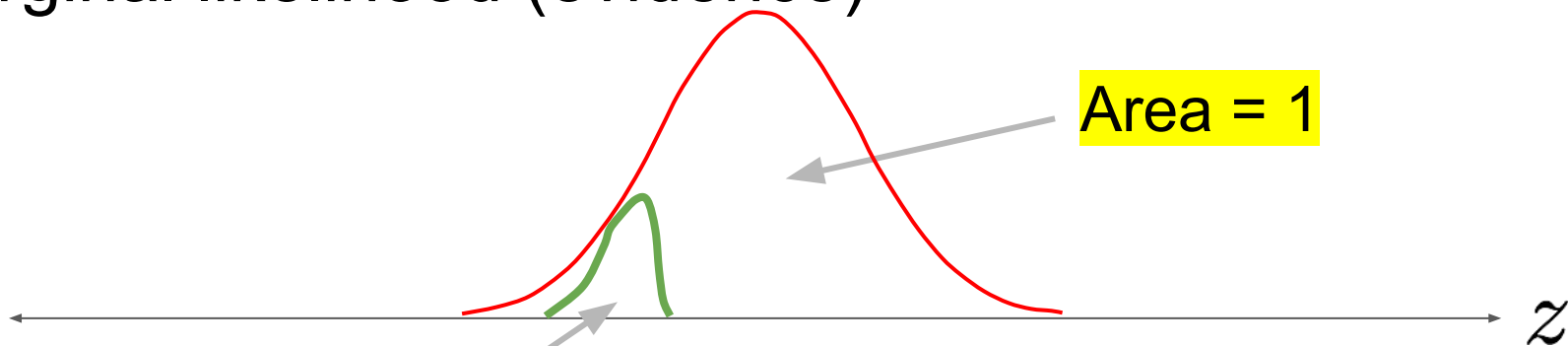
Joint density (with  $x$  observed)



1-minute exercise:  
what is the area?



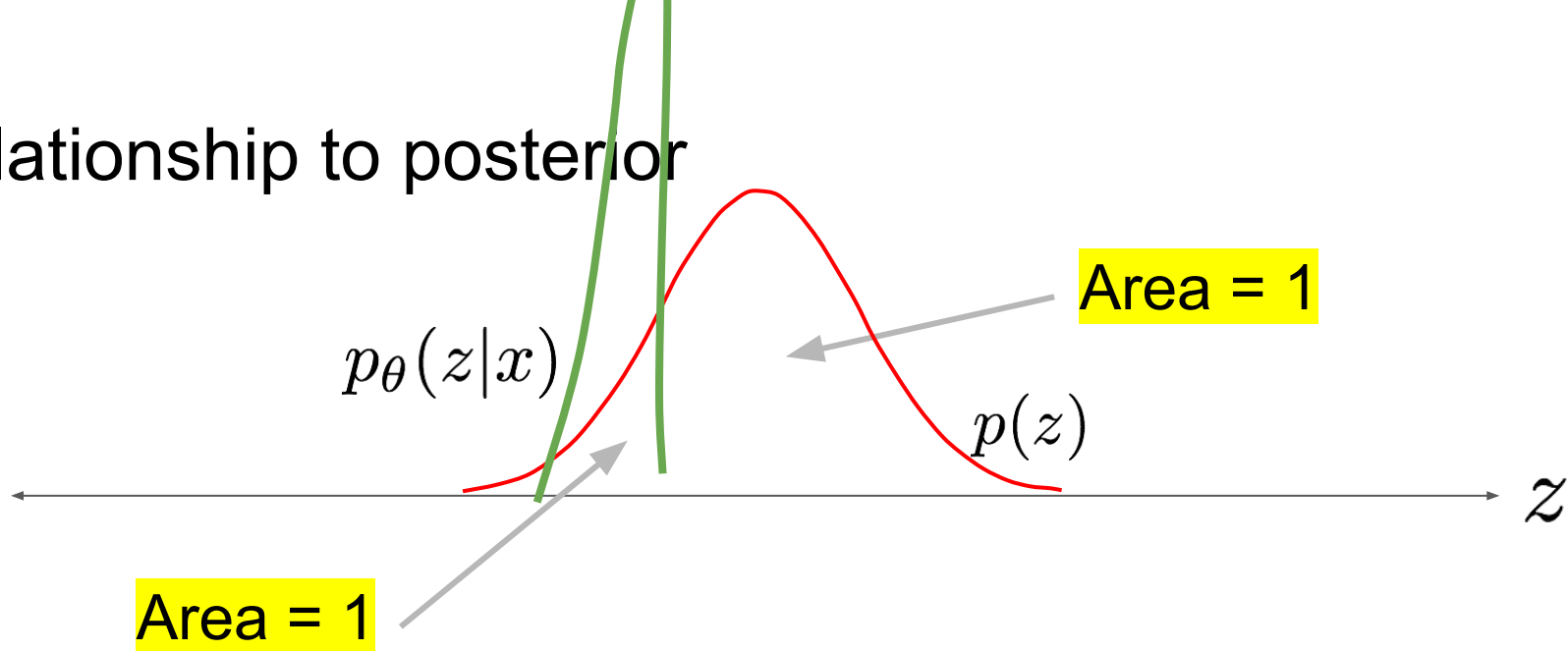
# Marginal likelihood (evidence)



Area = ?

$$\begin{aligned}\text{area} &= \int p_\theta(x|z) p(z) \, dz \\ &= \int p_\theta(x, z) \, dz \\ &= p_\theta(x)\end{aligned}$$

# Relationship to posterior



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

Dividing by the marginal likelihood (evidence) scales the area back to 1...

Evidence, for all data points

$$X \equiv x^{(1)}, x^{(2)} \dots, x^{(N)}$$

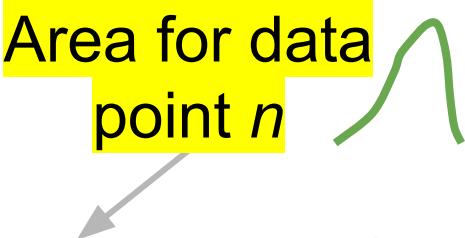
$$p_{\theta}(X) = \prod_{n=1}^N p_{\theta}(x^{(n)})$$

Area for data  
point  $n$



Evidence, for all data points

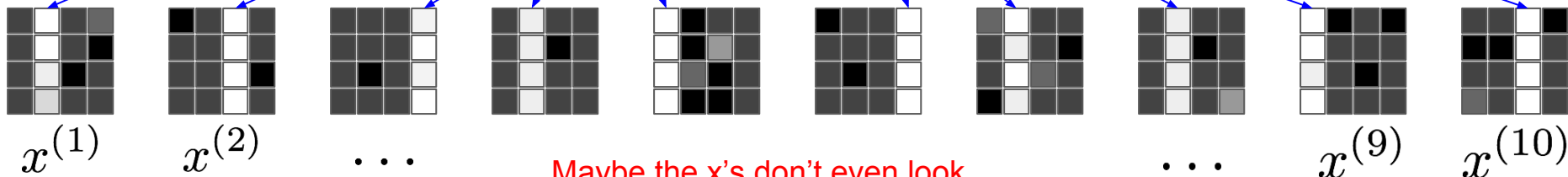
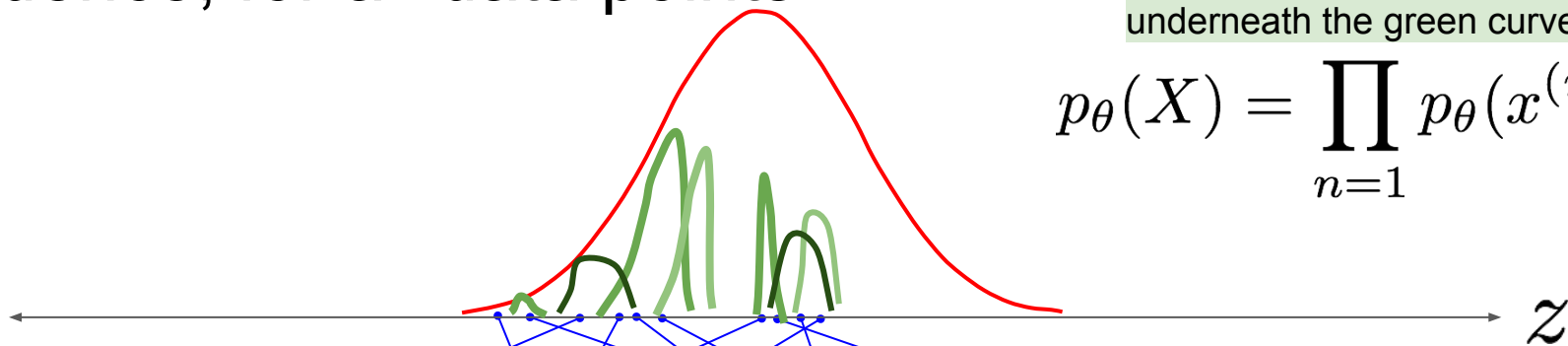
$$X \equiv x^{(1)}, x^{(2)} \dots, x^{(N)}$$


$$\log p_{\theta}(X) = \sum_{n=1}^N \log p_{\theta}(x^{(n)})$$


# Evidence, for all data points

The product of the areas underneath the green curves

$$p_{\theta}(X) = \prod_{n=1} p_{\theta}(x^{(n)})$$

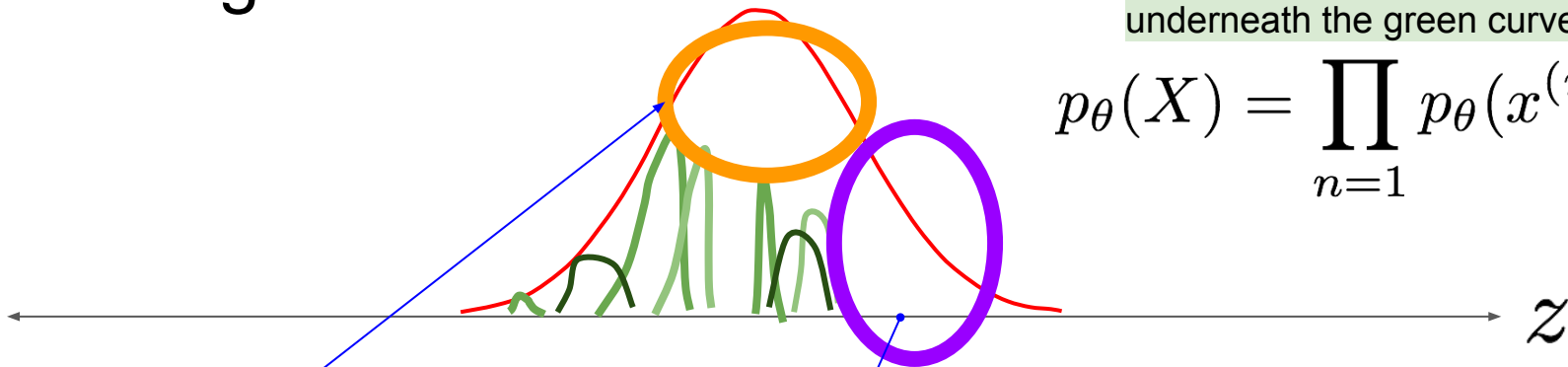


Maybe the  $x$ 's don't even look that nice, when sampled with 

# Maximizing the evidence

The product of the areas underneath the green curves

$$p_{\theta}(X) = \prod_{n=1} p_{\theta}(x^{(n)})$$



By changing  $\theta$  we can make the evidence for these data points bigger...

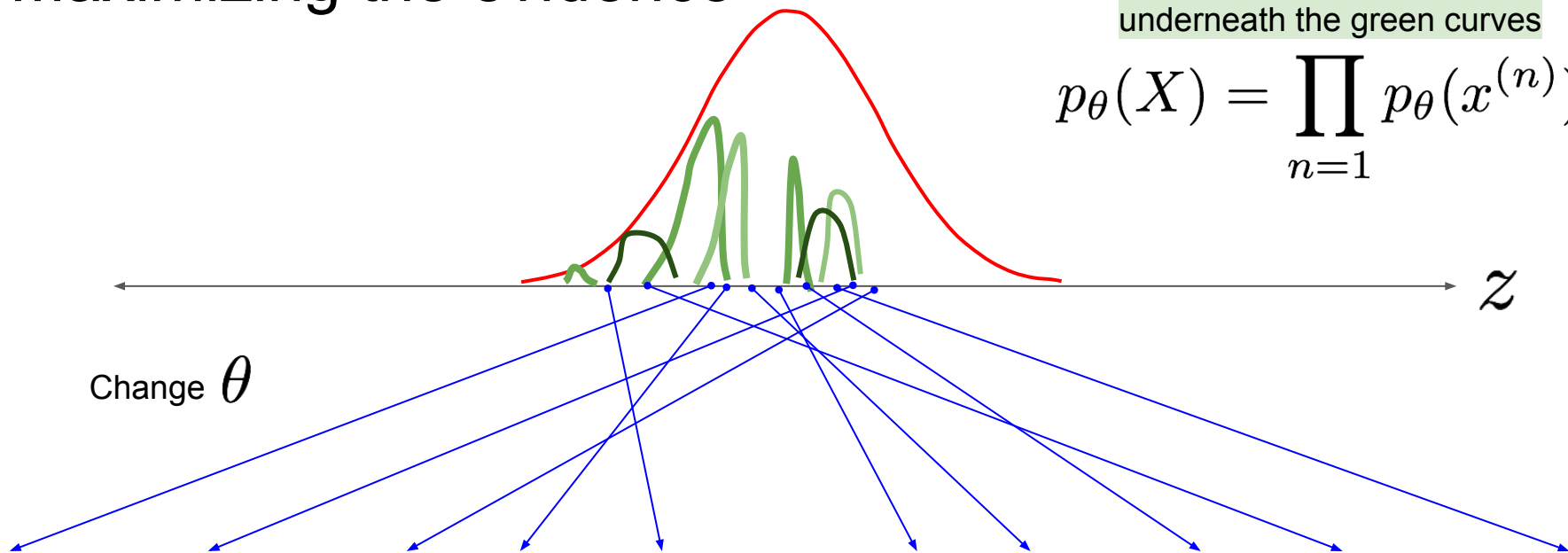
These  $z$ 's don't generate images like the ones in the data set...

(With this  $\theta$ , the prior doesn't capture the data manifold well)

# Maximizing the evidence

The product of the areas underneath the green curves

$$p_{\theta}(X) = \prod_{n=1} p_{\theta}(x^{(n)})$$

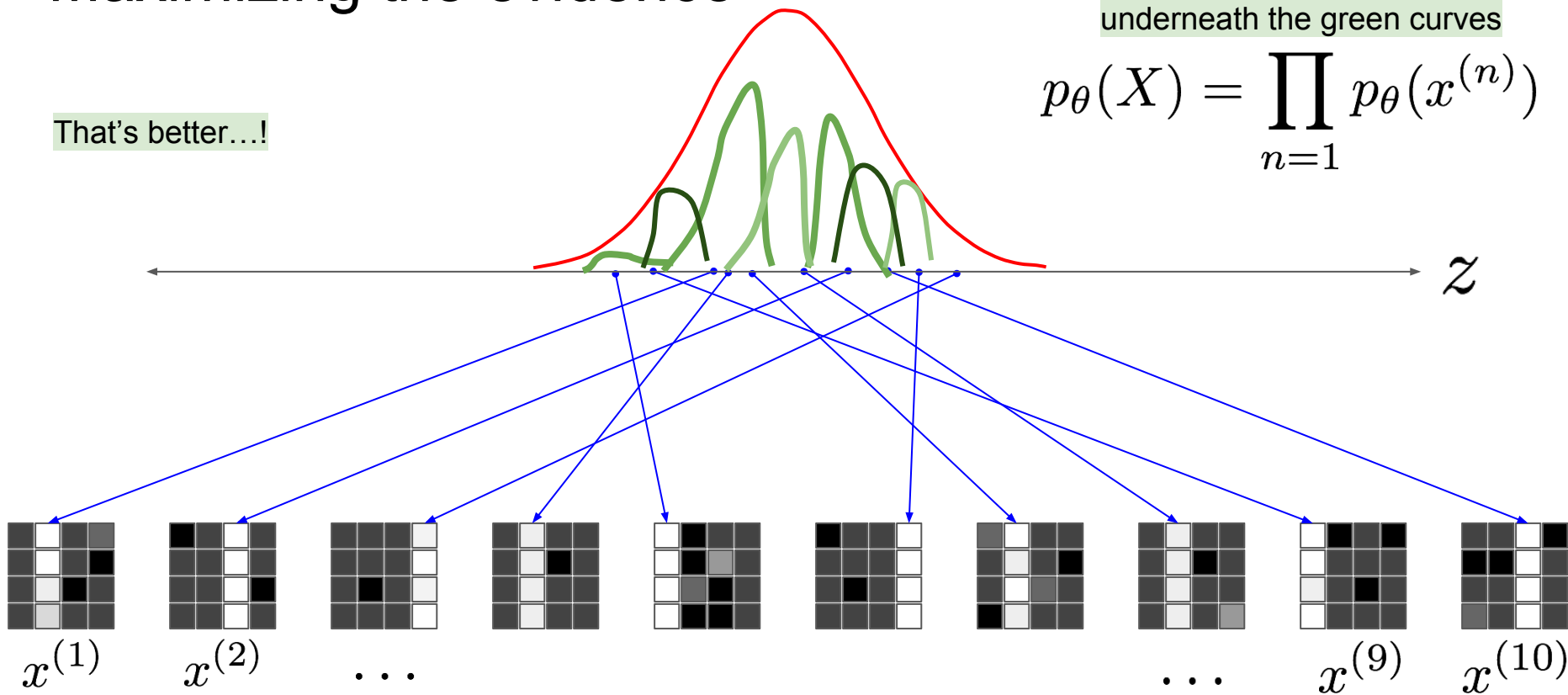


# Maximizing the evidence

That's better...!

The product of the areas underneath the green curves

$$p_{\theta}(X) = \prod_{n=1} p_{\theta}(x^{(n)})$$

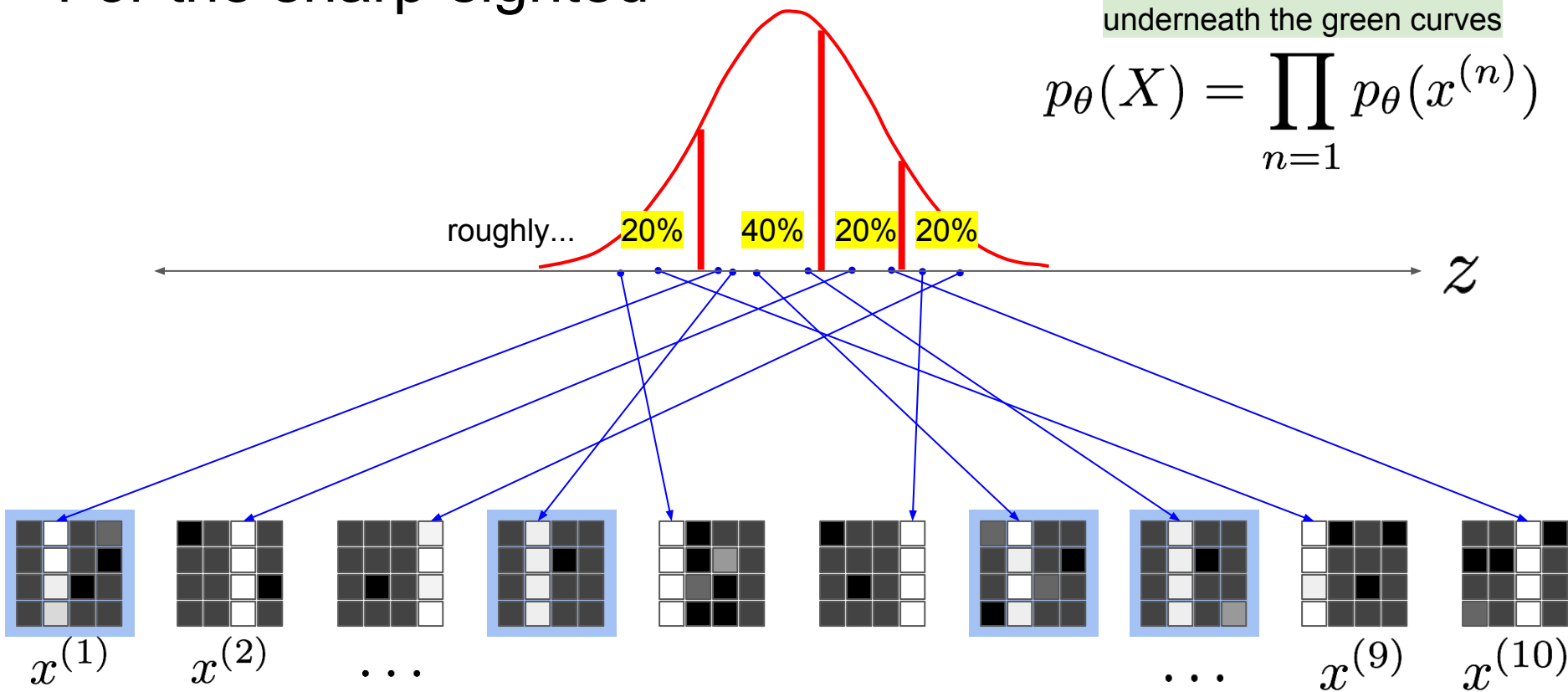




# For the sharp-sighted

The product of the areas underneath the green curves

$$p_{\theta}(X) = \prod_{n=1} p_{\theta}(x^{(n)})$$




# Learning

We want to maximize

$$\begin{aligned} & \max_{\theta} \left[ \log p_{\theta}(X) \right] \\ &= \max_{\theta} \left[ \sum_{n=1}^N \log p_{\theta}(x^{(n)}) \right] \end{aligned}$$

Area for data point  $n$



except that we cannot write down an analytically tractable expression for the area.

Strategies: Stochastic (Monte Carlo samples + gradients) or deterministic (approximate inference). We'll follow the "deterministic" path next...

# Approximate inference

We want to use this quantity for “learning”, but cannot compute it in an analytically tractable way:

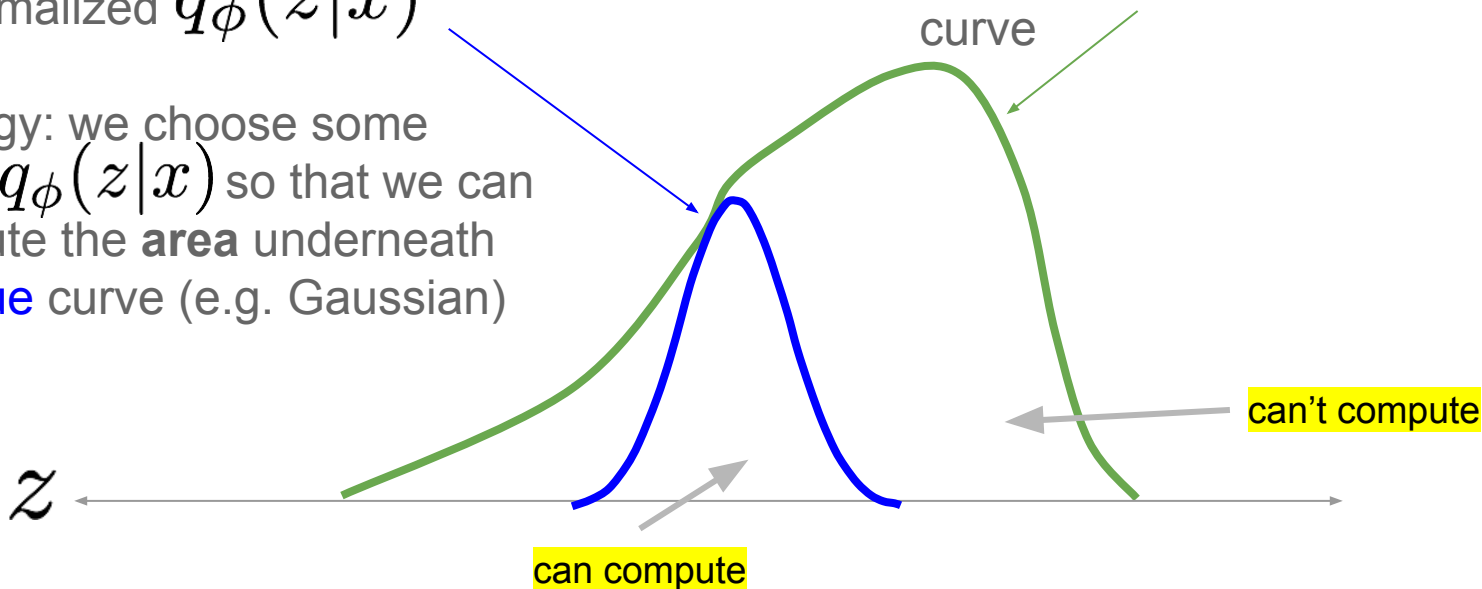
$$\begin{aligned}\log p_{\theta}(x) &= \log \int p_{\theta}(x|z) p(z) \, dz \\ &= \log \int p_{\theta}(x, z) \, dz\end{aligned}$$

# “Variational lower bound”

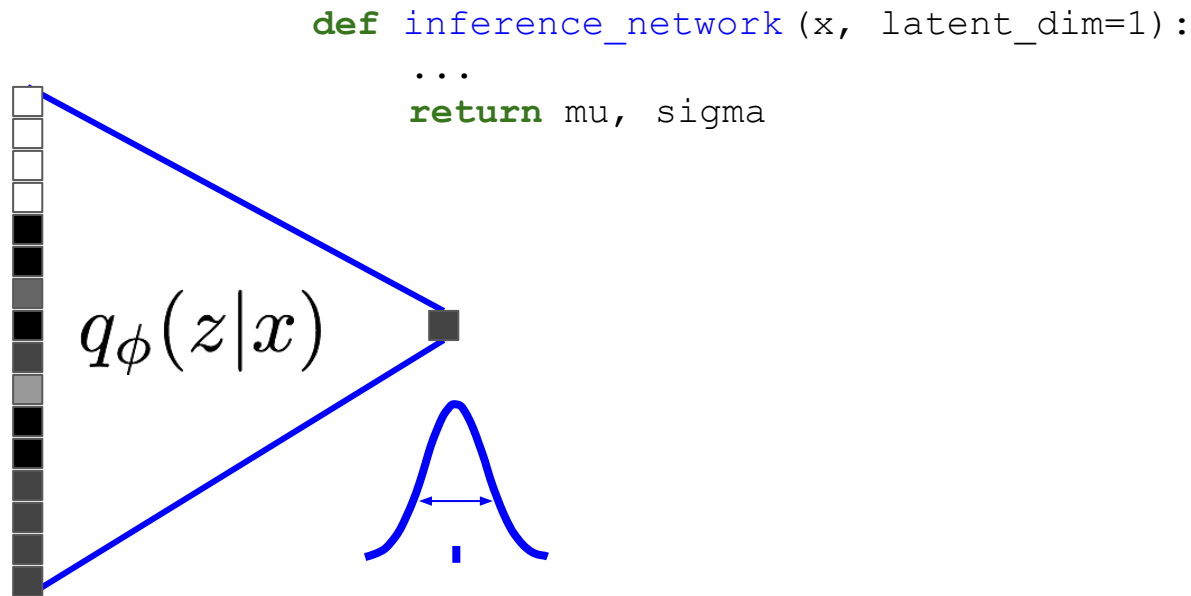
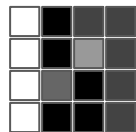
Unnormalized  $q_{\phi}(z|x)$

Strategy: we choose some other  $q_{\phi}(z|x)$  so that we can compute the **area** underneath the **blue** curve (e.g. Gaussian)

Unnormalized  $p_{\theta}(z|x)$   
We cannot (tractably) compute the **area** underneath the **green** curve

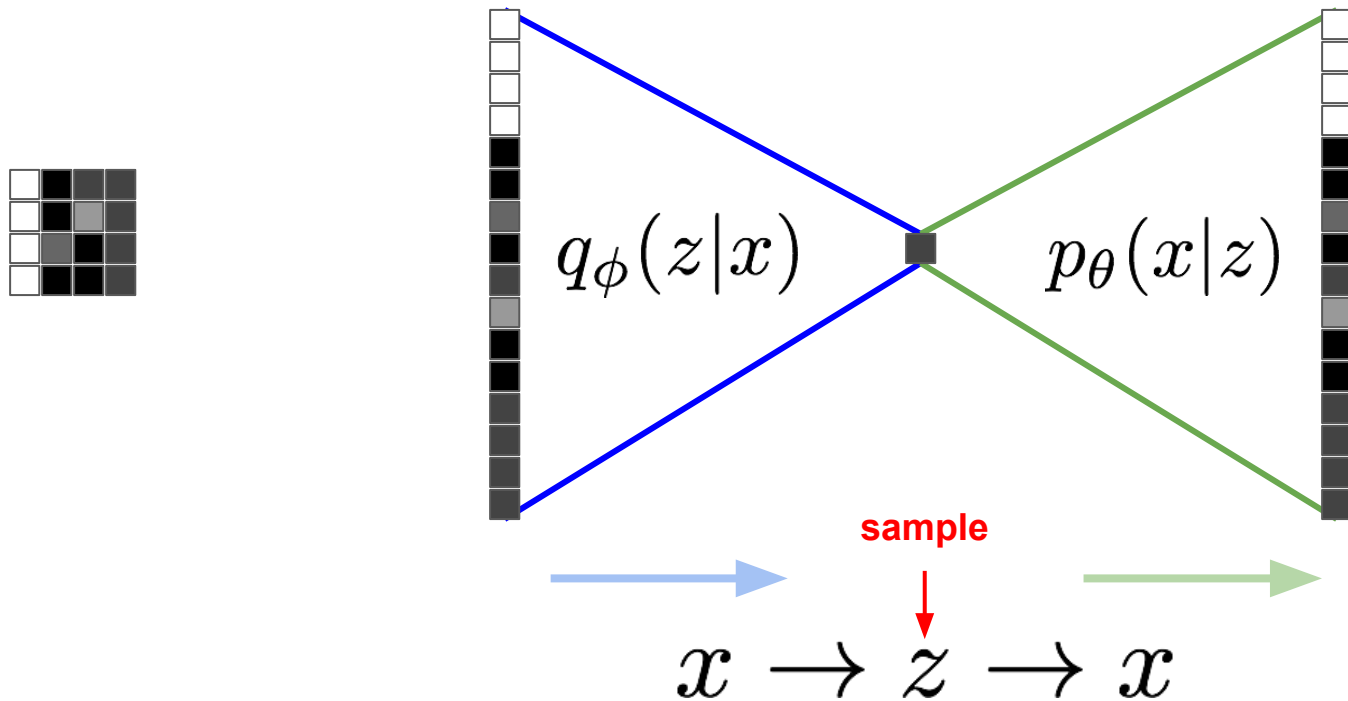


# Encoder



$$x \rightarrow \mu_{\phi}(x), \sigma_{\phi}^2(x)$$

# Encoder decoder

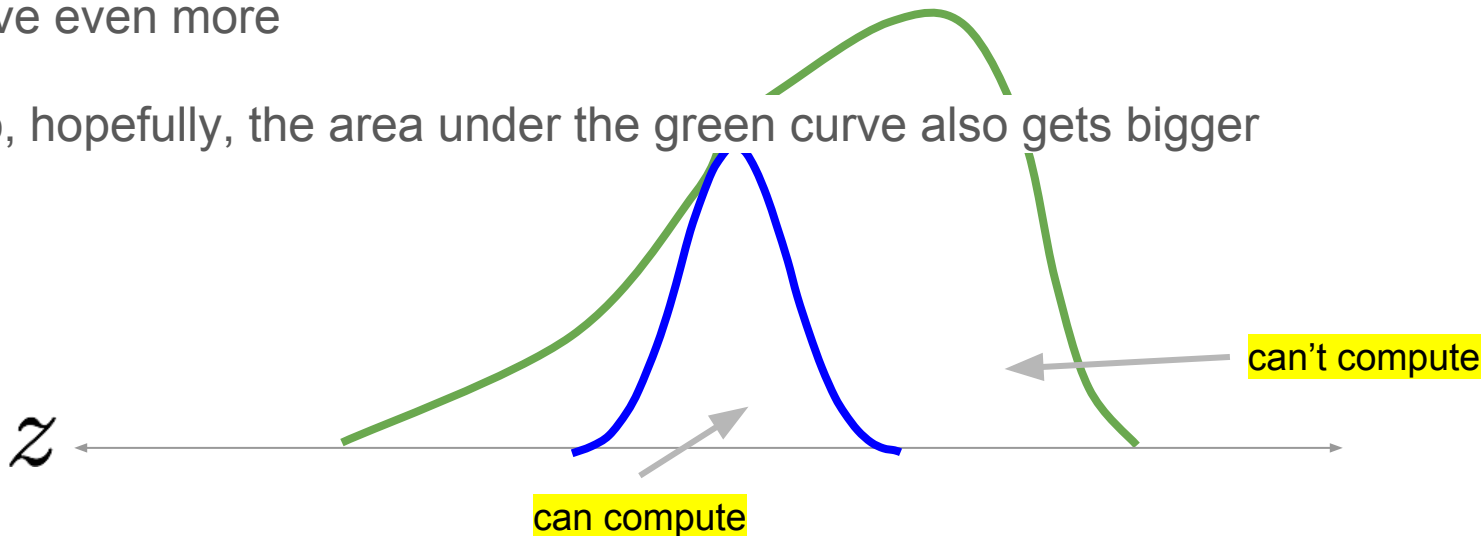


# Strategy

Change  $\phi$  to inflate the area under the blue curve. We can do that!

Change  $\theta$  to change the green curve, so that we can inflate the area under the blue curve even more

...and so, hopefully, the area under the green curve also gets bigger



Whhaaaatttt?

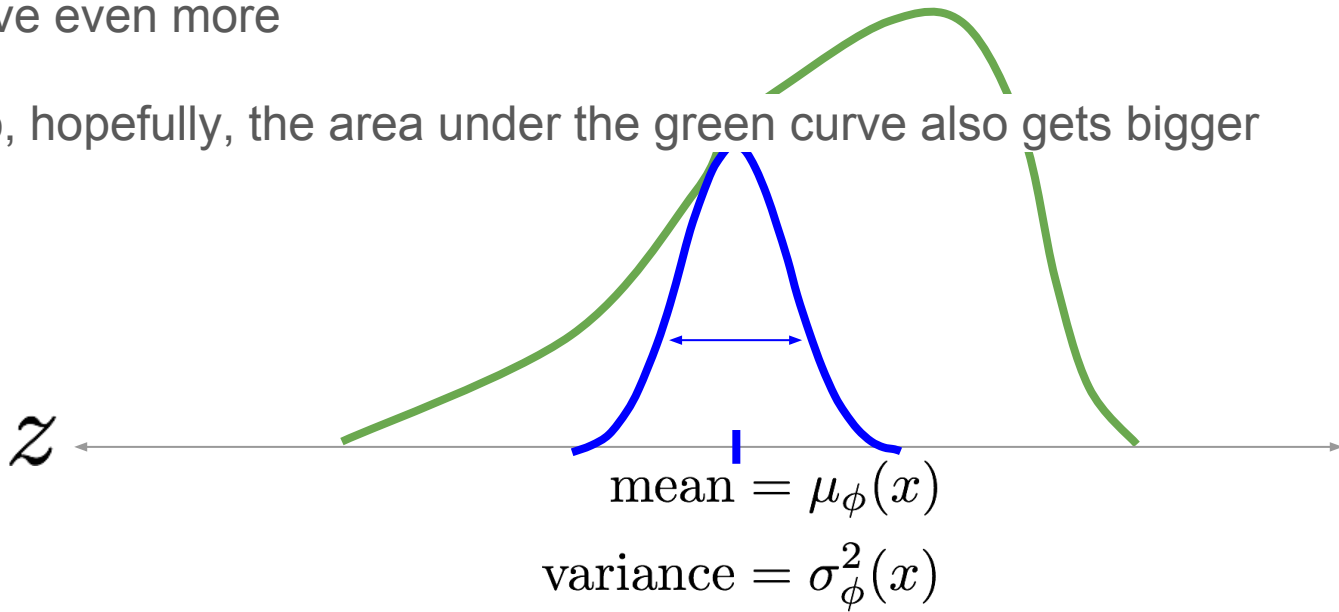


# Strategy

Change  $\phi$  to inflate the area under the blue curve. We can do that!

Change  $\theta$  to change the green curve, so that we can inflate the area under the blue curve even more

...and so, hopefully, the area under the green curve also gets bigger



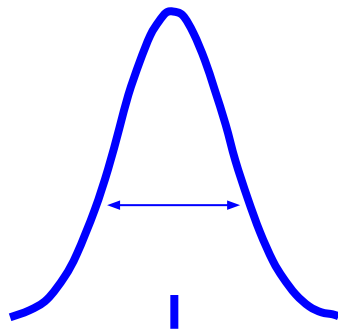
## 3-minute exercise

Create and draw  $q_{\phi}(z|x) = \mathcal{N}\left(z; \mu_{\phi}(x), \sigma_{\phi}^2(x)\right)$  as a function.

It could be a multi-layer perceptron (MLP) that takes 16-dimensional  $\mathbf{x}$ , and produces two 1-dimensional quantities,

$$\text{mean} = \mu_{\phi}(x)$$

$$\text{variance} = \sigma_{\phi}^2(x)$$

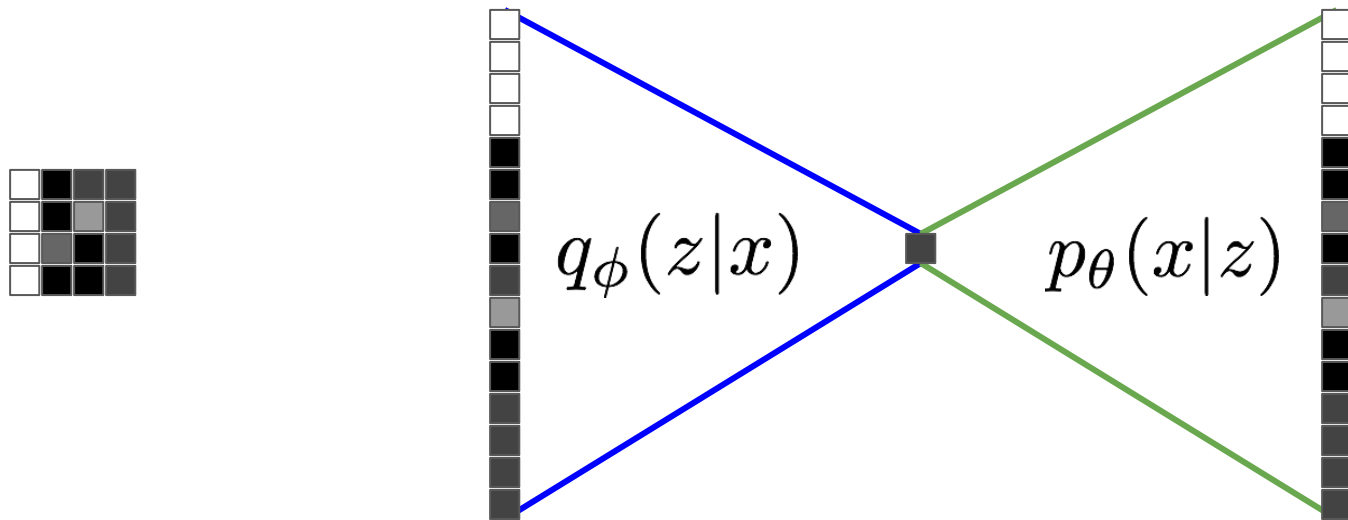


scratch space

What are your parameters  $\phi$ ?

# Objective function discussion

maximize (for all data points)...



$$\mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] - \text{KL}(q_\phi(z|x) \parallel p(z))$$

$\text{KL} \geq 0$

# Evidence lower bound (ELBO) for one data point

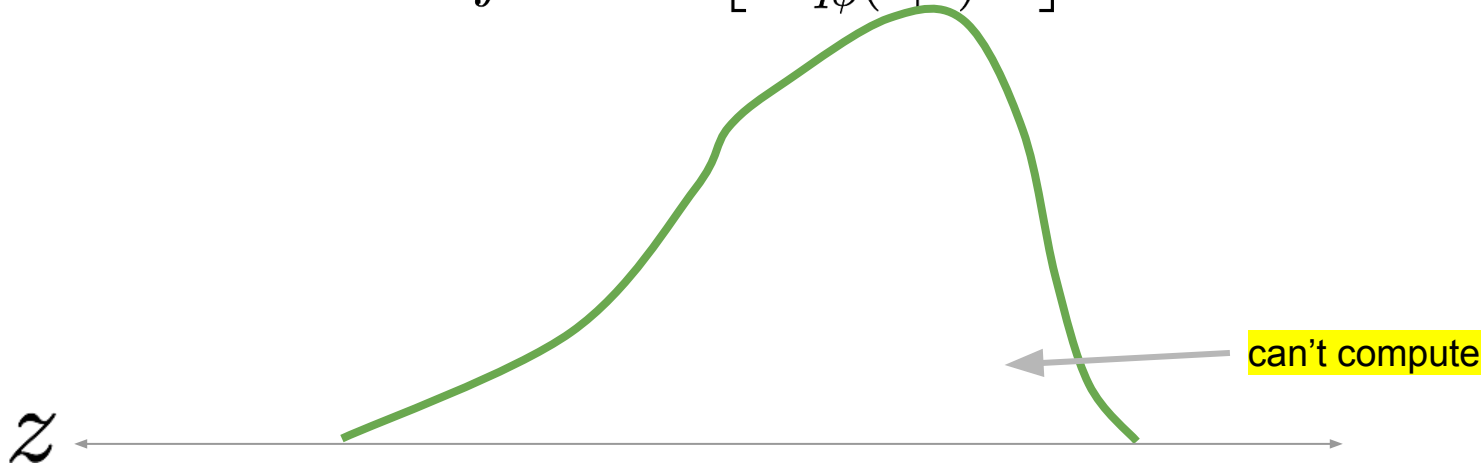
$$\begin{aligned}\log p_{\theta}(x) &= \log \int p_{\theta}(x|z) p(z) \mathrm{d}z \\&= \log \int q_{\phi}(z|x) \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] \mathrm{d}z \\&\geq \int q_{\phi}(z|x) \log \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] \mathrm{d}z \quad [\text{Jensen}] \\&= \int q_{\phi}(z|x) \log p_{\theta}(x|z) \mathrm{d}z - \int q_{\phi}(z|x) \log \left[ \frac{q_{\phi}(z|x)}{p(z)} \right] \mathrm{d}z \\&= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log p_{\theta}(x|z) \right] - \text{KL}(q_{\phi}(z|x) \parallel p(z)) \\&\equiv \mathcal{L}(x; \theta, \phi)\end{aligned}$$

ELBO



# Evidence lower bound (ELBO) for one data point

$$\begin{aligned}\log p_{\theta}(x) &= \log \int p_{\theta}(x|z) p(z) \mathrm{d}z \\ &= \log \int q_{\phi}(z|x) \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] \mathrm{d}z\end{aligned}$$



# Evidence lower bound (ELBO) for one data point

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

concave function

$$= \log \int \underline{q_{\phi}(z|x)} \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] \underline{dz}$$
$$\geq \int \underline{q_{\phi}(z|x)} \log \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] \underline{dz} \quad [\text{Jensen}]$$

expectation      concave function

$$\log \mathbb{E}_{q_{\phi}(z|x)} \left[ f(z) \right] \geq \mathbb{E}_{q_{\phi}(z|x)} \left[ \log f(z) \right]$$

# 3-minute exercise

Jensen's inequality

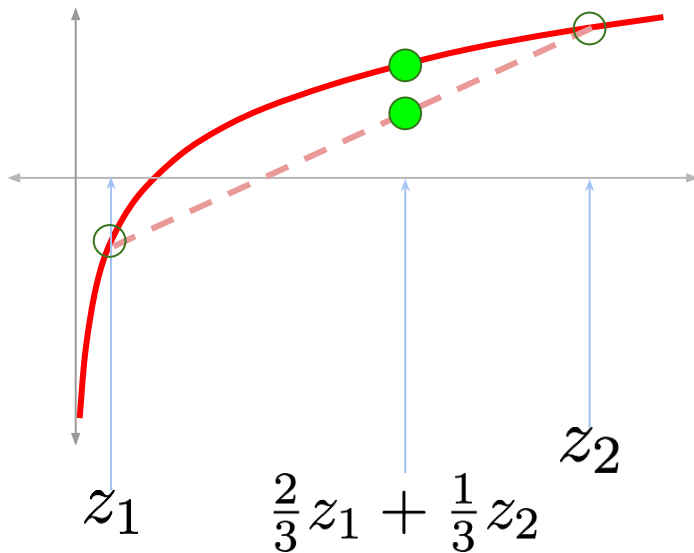
Draw  $\log(\dots)$  as a function, and convince yourself that

$$\log \left( \underline{\frac{2}{3}} z_1 + \underline{\frac{1}{3}} z_2 \right) \geq \underline{\frac{2}{3}} \log(z_1) + \underline{\frac{1}{3}} \log(z_2)$$

is always true for any (nonnegative) setting of  $z_1$  and  $z_2$ .

# Logarithm (concave)

$$\log\left(\frac{2}{3}z_1 + \frac{1}{3}z_2\right) \geq \frac{2}{3}\log(z_1) + \frac{1}{3}\log(z_2)$$





# Evidence lower bound (ELBO) for one data point

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

$$= \log \int q_{\phi}(z|x) \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] dz$$

Kullback-Leibler  
divergence between  
two Gaussian  
distributions (here).  
Can compute in closed  
form

## Reconstruction

Expected log likelihood.  
Cannot compute in  
closed form, and will  
have to get a Monte  
Carlo estimate  
(with SGD)

$$\geq \int q_{\phi}(z|x) \log \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] dz$$

$$= \int q_{\phi}(z|x) \log p_{\theta}(x|z) dz - \int q_{\phi}(z|x) \log \left[ \frac{q_{\phi}(z|x)}{p(z)} \right] dz$$

$$= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log p_{\theta}(x|z) \right] - \text{KL}(q_{\phi}(z|x) \parallel p(z))$$

ELBO

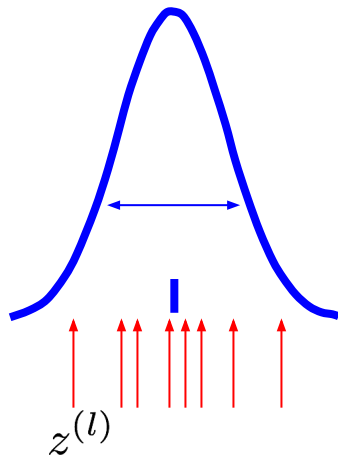


$$\equiv \mathcal{L}(x; \theta, \phi)$$

# Expected log likelihood

We can estimate the expected log likelihood with a Monte Carlo estimate:


Draw L samples  $z^{(l)} \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x)) \dots$



# Expected log likelihood

We can estimate the expected log likelihood with a Monte Carlo estimate:

Draw  $L$  samples  $z^{(l)} \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))$  and use them to estimate the average:

$$\begin{aligned}\mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] &= \mathbb{E}_{z \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))} \left[ \log p_\theta(x|z) \right] \\ &\approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(x|z^{(l)})\end{aligned}$$


# Expected log likelihood

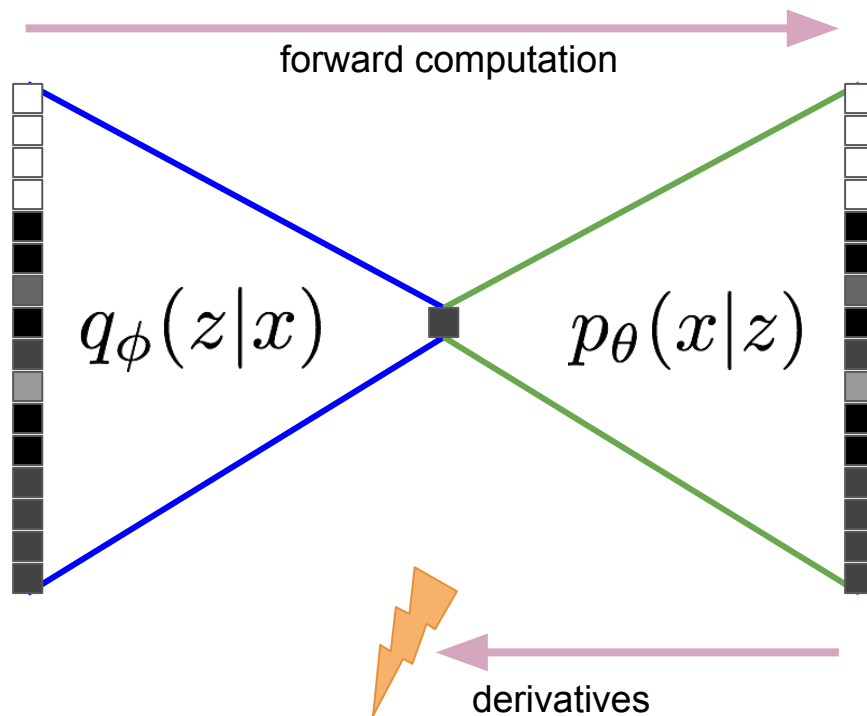
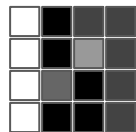
We can estimate the expected log likelihood with a Monte Carlo estimate:

Draw  $L$  samples  $z^{(l)} \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))$  and use them to estimate the average:

$$\begin{aligned}\mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] &= \mathbb{E}_{z \sim \mathcal{N}(z; \mu_\phi(x), \sigma_\phi^2(x))} \left[ \log p_\theta(x|z) \right] \\ &\approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(x|z^{(l)})\end{aligned}$$

Using samples in *this* way removes  $\phi$  from part of the objective function, and even though we can evaluate it, we can't take derivatives / get the gradients!

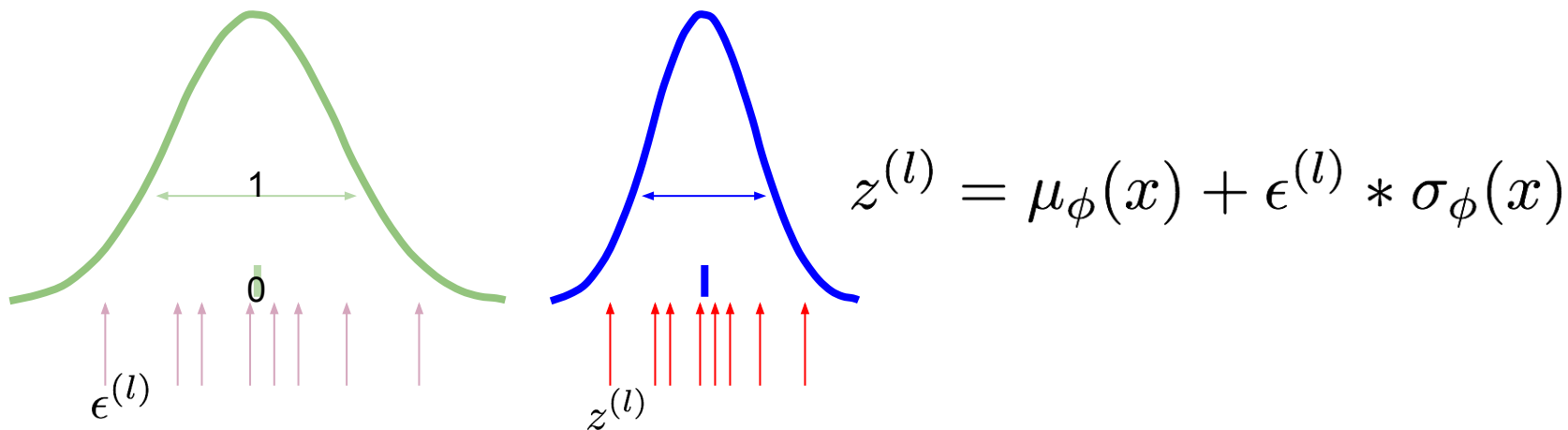
# Naive sampling



# Expected log likelihood: reparameterization trick

We can estimate the expected log likelihood with a Monte Carlo estimate:


Draw  $L$  samples  $\epsilon^{(l)} \sim \mathcal{N}(\epsilon; 0, 1)$  and transform them!



# Expected log likelihood

We can estimate the expected log likelihood with a Monte Carlo estimate:

Draw  $L$  samples  $\epsilon^{(l)} \sim \mathcal{N}(\epsilon; 0, 1)$  and use them to estimate the average:

$$\begin{aligned}\mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] &= \mathbb{E}_{\epsilon \sim \mathcal{N}(\epsilon; 0, 1)} \left[ \log p_\theta(x \mid z = \mu_\phi(x) + \epsilon * \sigma_\phi(x)) \right] \\ &\approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(x \mid z^{(l)} = \mu_\phi(x) + \epsilon^{(l)} * \sigma_\phi(x))\end{aligned}$$


# Expected log likelihood

We can estimate the expected log likelihood with a Monte Carlo estimate:

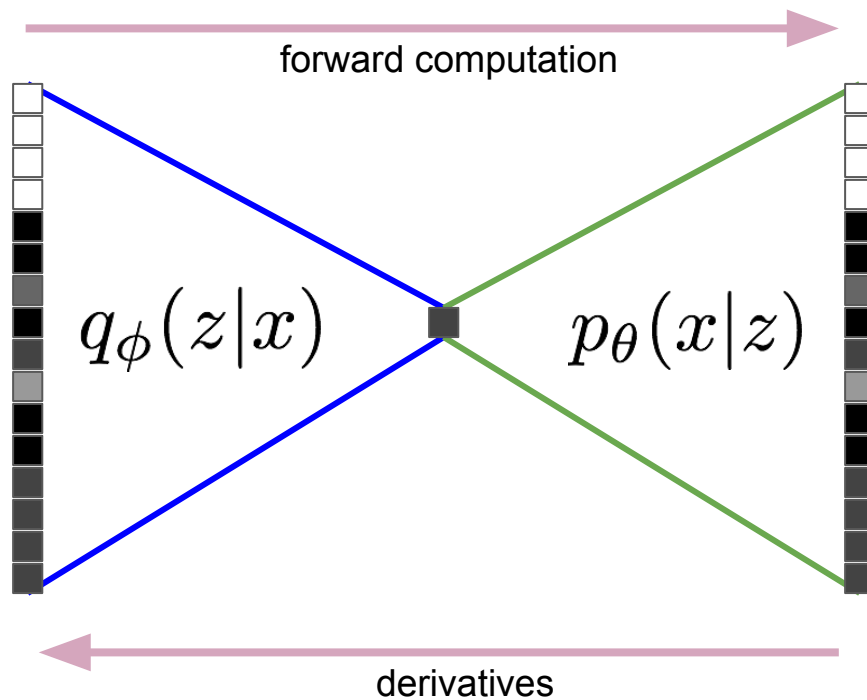
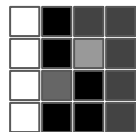
Draw  $L$  samples  $\epsilon^{(l)} \sim \mathcal{N}(\epsilon; 0, 1)$  and use them to estimate the average:

$$\begin{aligned}\mathbb{E}_{q_\phi(z|x)} \left[ \log p_\theta(x|z) \right] &= \mathbb{E}_{\epsilon \sim \mathcal{N}(\epsilon; 0, 1)} \left[ \log p_\theta(x \mid z = \mu_\phi(x) + \epsilon * \sigma_\phi(x)) \right] \\ &\approx \frac{1}{L} \sum_{l=1}^L \log p_\theta(x \mid z^{(l)} = \mu_\phi(x) + \epsilon^{(l)} * \sigma_\phi(x))\end{aligned}$$

The noise is introduced “from outside” the computation graph, and we can evaluate the objective function **and** take derivatives / get the gradients!



# Reparameterization trick

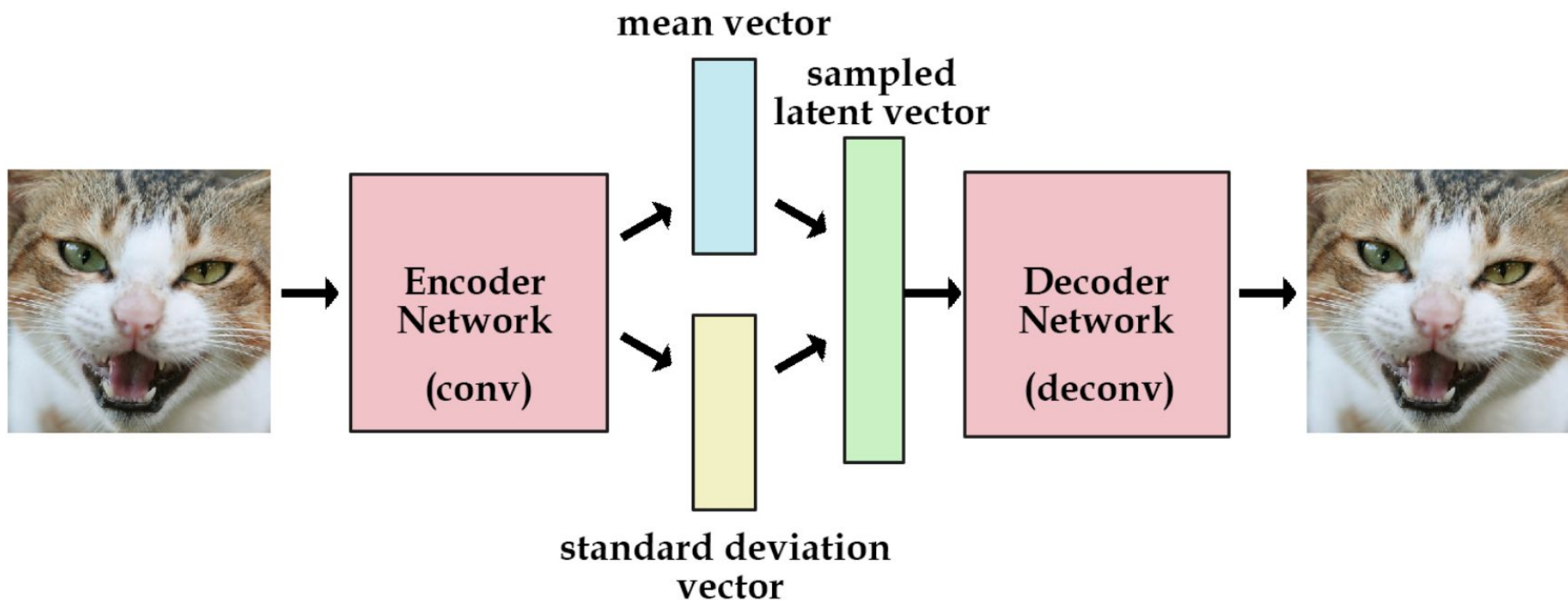


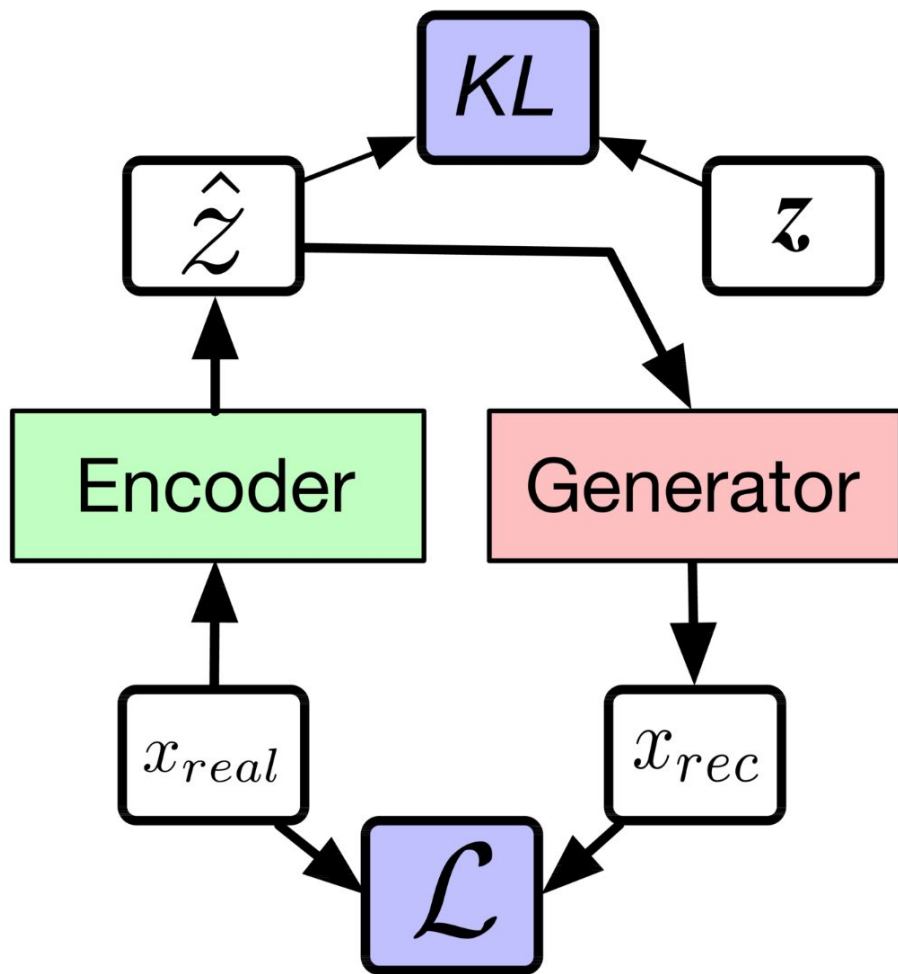
# ELBO for full data set

You now have all the tools to estimate the ELBO for a whole data set,

$$\mathcal{L}(X; \theta, \phi) = \sum_{n=1}^N \left\{ \mathbb{E}_{q_{\phi}(z^{(n)}|x^{(n)})} \left[ \log p_{\theta}(x^{(n)}|z^{(n)}) \right] - \text{KL}(q_{\phi}(z^{(n)}|x^{(n)}) \parallel p(z^{(n)})) \right\}$$

take mini-batch subsamples, and use stochastic gradient ascent to maximize it.

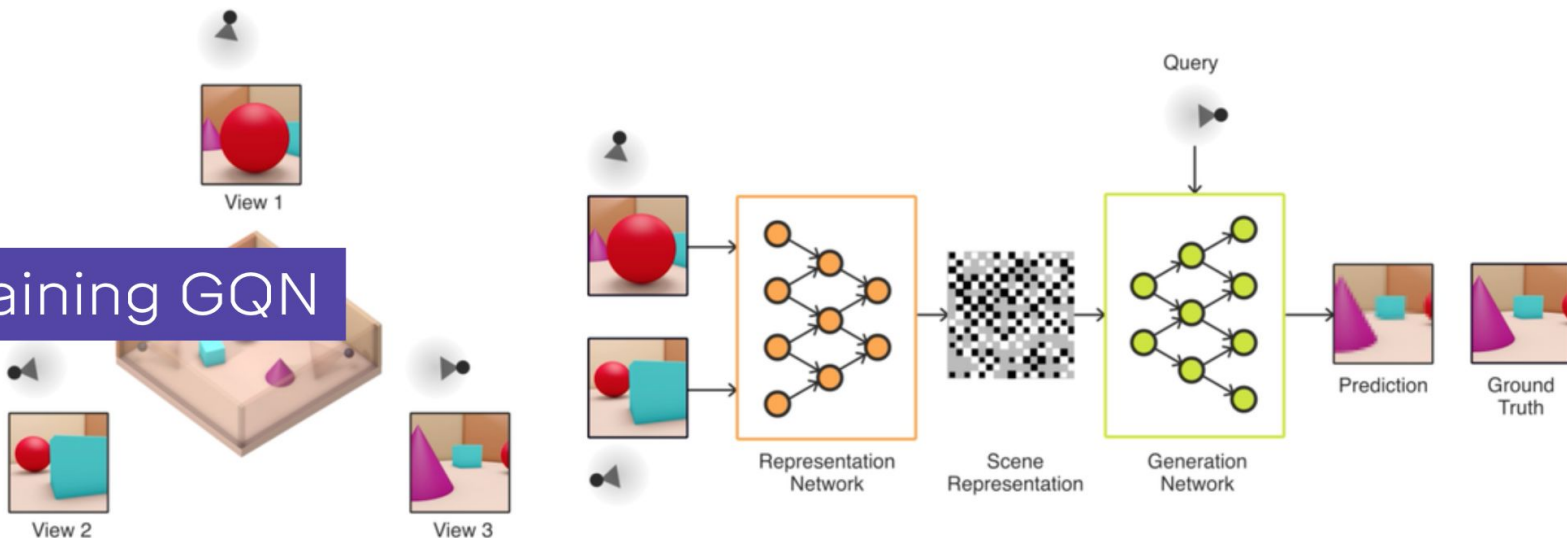






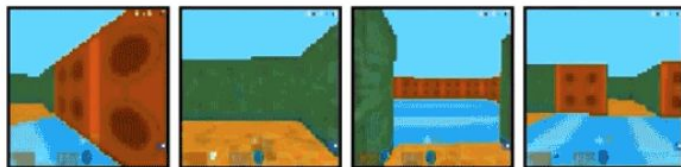
# Neural scene representation and rendering

## ► Explaining GQN



# Neural scene representation and rendering

observations



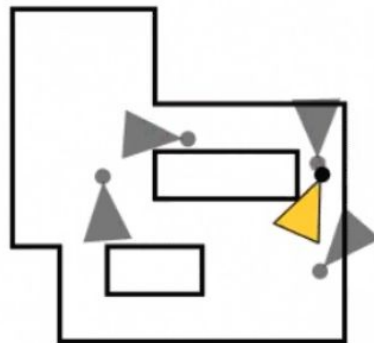
ground truth



neural rendering

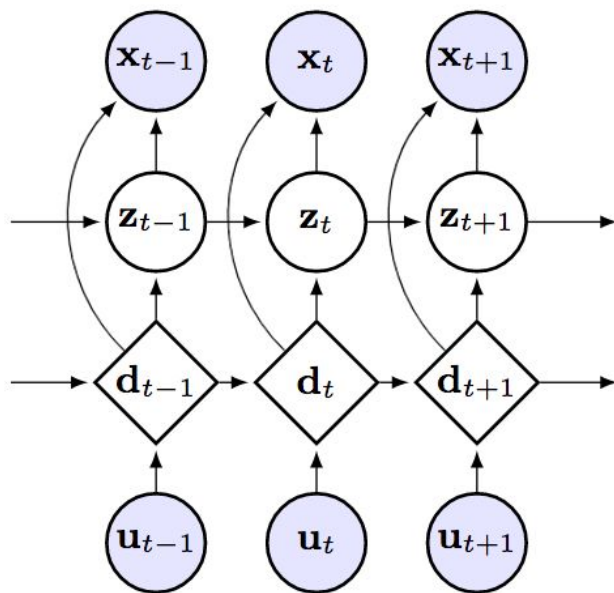


map



# Sequential Neural Models with Stochastic Layers

## Stochastic Recurrent Neural Networks



(a) Generative model  $p_\theta$

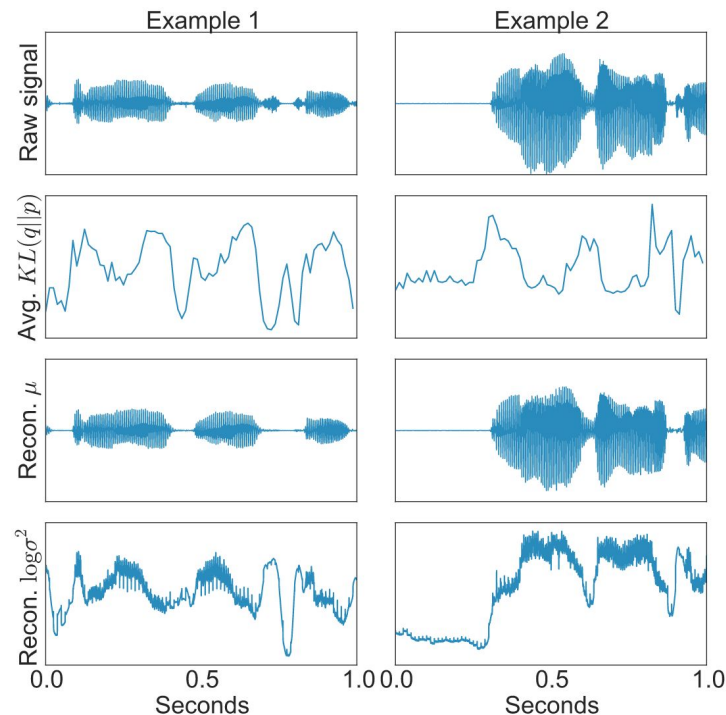
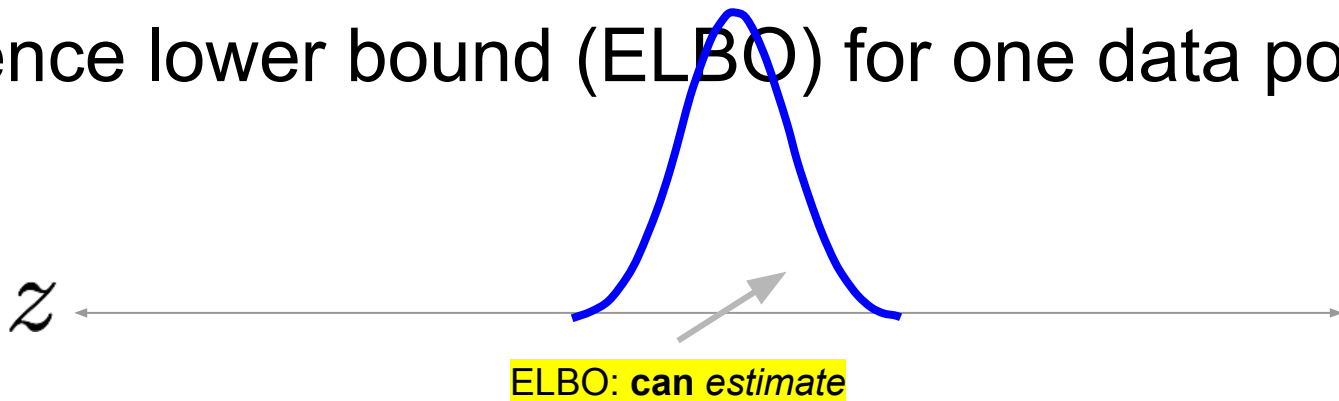


Figure 3: Visualization of the average KL term and reconstructions of the output mean and log-variance for two examples from the Blizzard test set.



The end

# Evidence lower bound (ELBO) for one data point



$$\begin{aligned} &\geq \int q_{\phi}(z|x) \log \left[ \frac{p_{\theta}(x|z) p(z)}{q_{\phi}(z|x)} \right] dz \quad [\text{Jensen}] \\ &= \int q_{\phi}(z|x) \log p_{\theta}(x|z) dz - \int q_{\phi}(z|x) \log \left[ \frac{q_{\phi}(z|x)}{p(z)} \right] dz \\ &= \mathbb{E}_{q_{\phi}(z|x)} \left[ \log p_{\theta}(x|z) \right] - \text{KL}(q_{\phi}(z|x) \parallel p(z)) \\ &\equiv \mathcal{L}(x; \theta, \phi) \end{aligned}$$