# Lecture 9: Explicit Generative Models
## Efstratios Gavves

# Lecture overview

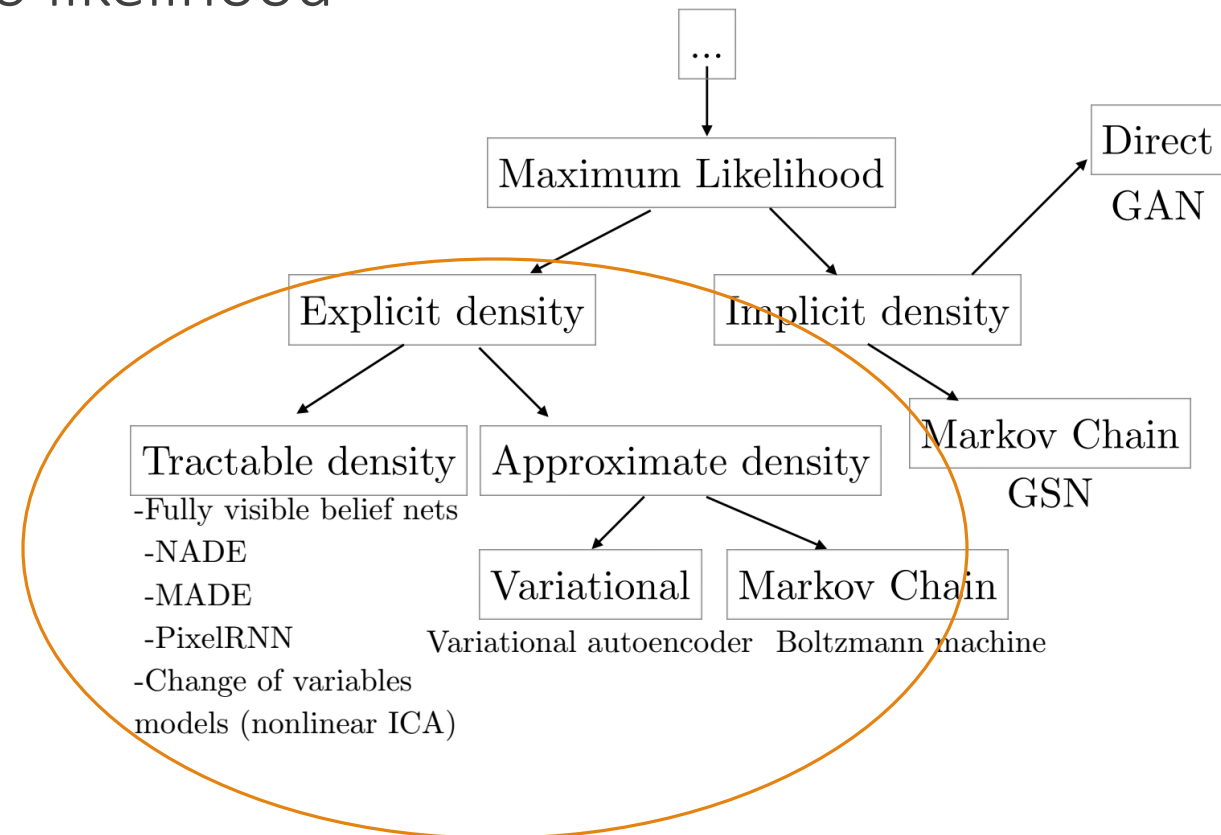o Gentle intro to Bayesian Modelling and Variational Inference

o Restricted Boltzmann Machines

o Deep Boltzmann Machines

o Deep Belief Network

o Contrastive Divergence

o Variational Autoencoders

o Normalizing Flows

# Explicit density models

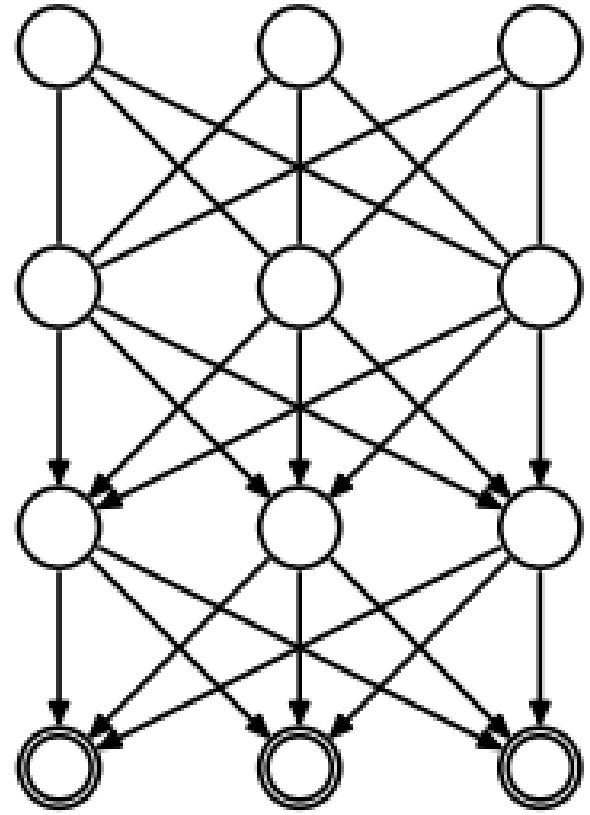o Plug in the model density function to likelihood
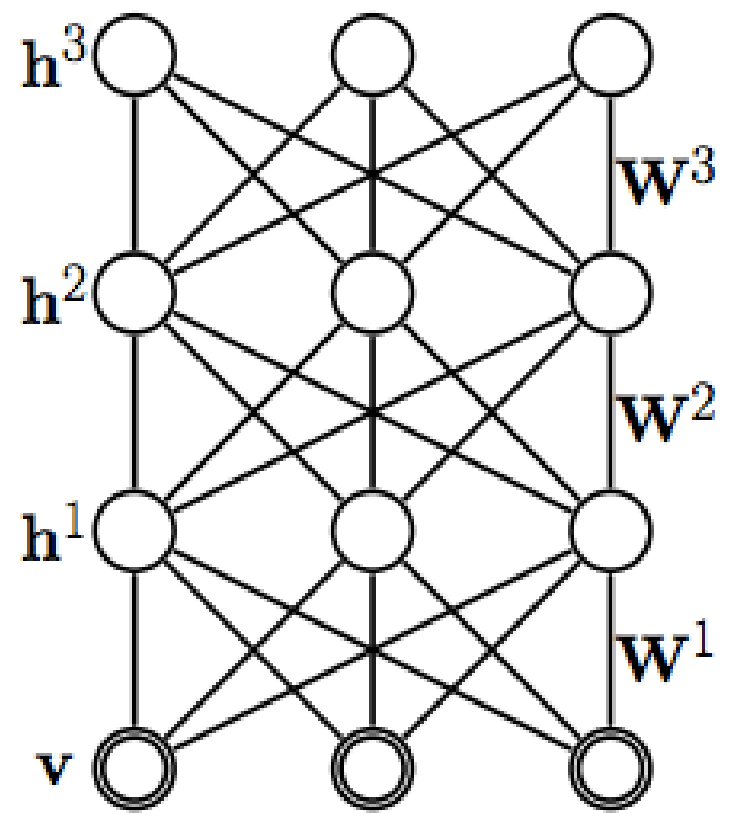
o Then maximize the likelihood

Problems

o Design complex enough model that meets data complexity

o At the same time, make sure model is computationally tractable

o More details in the next lecture

Bayesian Modelling
Variational Inference



**Deep Belief Network**

**Deep Boltzmann Machine**

$h^3$
$h^2$
$h^1$
$v$

$W^3$
$W^2$
$W^1$

# How to define a generative model?

o We can define an explicit density function over all possible relations $\psi_c$ between the input variables $x_c$

$$p(x) = \prod_c \psi_c(x_c)$$

o Quite inefficient → think of all possible relations (not just pairwise) between $256 \times 256 = 65K$ input variables

o Solution: Define an energy function to model the relations between the inputs variables

# Restricted Boltzmann Machines

o Boltzmann (or Gibbs) distribution defined over a free energy function $E(x)$

$$p(x) = \frac{1}{Z}\exp(-E(x))$$

o $Z$ is the normalization factor that makes sure $\int_x p(x)\, dx = 1$

◦ Very expensive to compute ➔ if $x = \{0, 1\}$ computing $Z$ requires $2^d$ computations

o Better restrict the model further to a bottleneck

$$E(x) = -x^T W h - b^T x - c^T h$$

# Why Boltzmann?

o In statistical mechanics and mathematics, a Boltzmann distribution (also called Gibbs distribution) is a probability distribution, probability measure, or frequency distribution of particles in a system over various possible states. The distribution is expressed in the form

$$F(state) \propto \exp(-\frac{E}{kT})$$

o $E$ is the state energy, $k$ is the Boltzmann constant, $T$ is the thermodynamic temperature

https://en.wikipedia.org/wiki/Boltzmann_distribution

# Restricted Boltzmann Machines

- $E(x) = -x^T W h - b^T x - c^T h$

- The $x^T W h$ models correlations between $x$ and the latent activations via the parameter matrix $W$

- The $b^T x, c^T h$ model the priors

- Restricted Boltzmann Machines (RBM) assume $x, h$ to be binary

# Restricted Boltzmann Machines

○ $E(x) = -x^T W h - b^T x - c^T h, \quad \theta = \{W, b, c\}$

○ The free energy function $F(x) = -\log \sum_h \exp(-E(x,h))$
defines a bipartite graph
with undirected connections
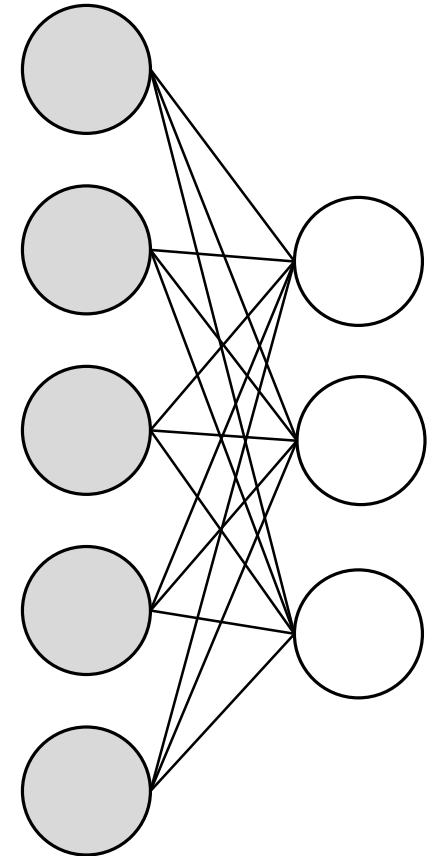◦ Information flows forward and backward

# Restricted Boltzmann Machines

○ The hidden units $h_j$ are independent to each other conditioned on the visible units

$$p(h|x) = \prod_j p(h_j|x, \theta)$$

○ The hidden units $x_i$ are independent to each other conditioned on the visible units

$$p(x|h) = \prod_i p(x_i|h, \theta)$$

# Training RBMs

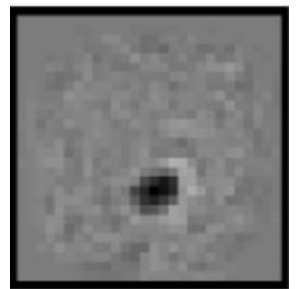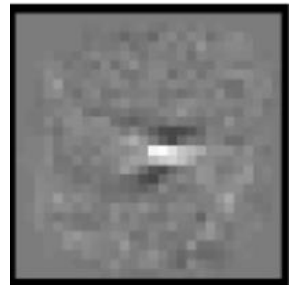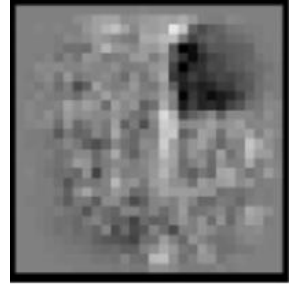o The conditional probabilities are defined as sigmoids

$$p(h_j|x,\theta) = \sigma(W_{.j}x + b_j)$$
$$p(x_i|h,\theta) = \sigma(W_{.i}x + c_i)$$

o Maximize log-likelihood

$$\mathcal{L}(\theta) = \frac{1}{N}\sum_n \log p(x_n|\theta)$$

o Let's take the gradients

$$\frac{\partial \log p(x_n|\theta)}{\partial \theta} = -\frac{\partial F(x_n)}{\partial \theta} - \frac{\partial \log Z}{\partial \theta}$$

$$= -\sum_h p(h|x_n,\theta)\frac{\partial E(x_n|h,\theta)}{\partial \theta} + \sum_{\tilde{x},h} p(\tilde{x},h|\theta)\frac{\partial E(\tilde{x},h|\theta)}{\partial \theta}$$



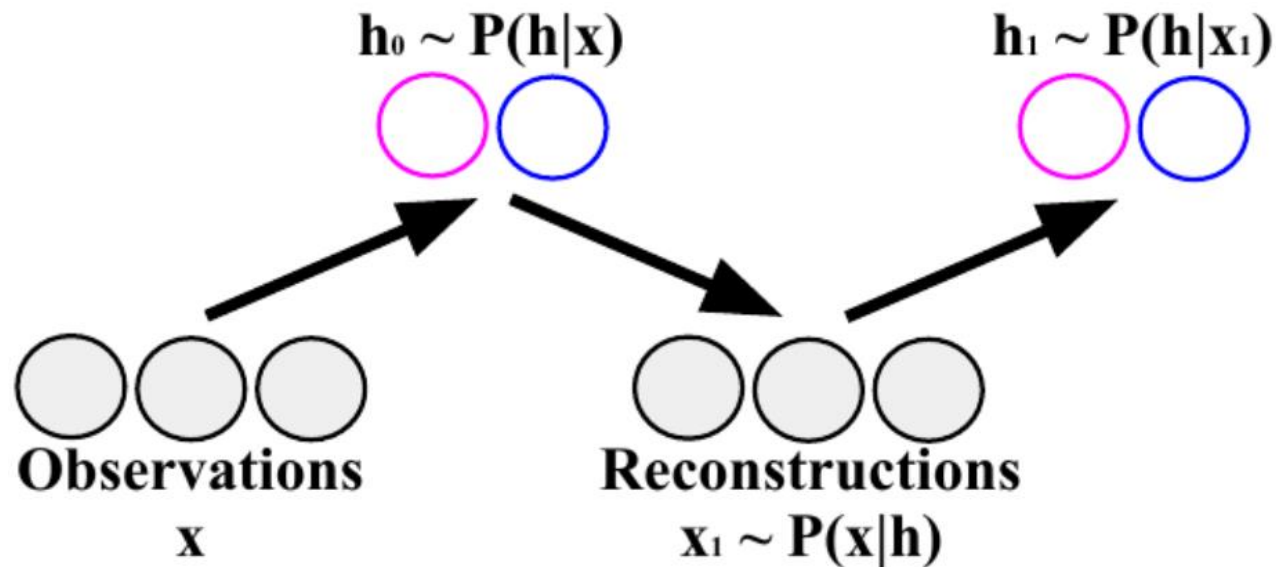Hidden unit (features)

# Training RBMs

- Let's take the gradients

$$\frac{\partial \log p(x_n|\theta)}{\partial \theta} = -\frac{\partial F(x_n)}{\partial \theta} - \frac{\partial \log Z}{\partial \theta}$$

$$= -\sum_h p(h|x_n, \theta)\frac{\partial E(x_n|h, \theta)}{\partial \theta} + \sum_{\tilde{x}, h} p(\tilde{x}, h|\theta)\frac{\partial E(\tilde{x}, h|\theta)}{\partial \theta}$$

- Easy because we just substitute in the definitions the $x_n$ and sum over $h$

- Hard because you need to sum over both $\tilde{x}, h$ which can be huge
  - It requires approximate inference, *e.g.,* MCMC

# Training RBMs with Contrastive Divergence

o Approximate the gradient with Contrastive Divergence

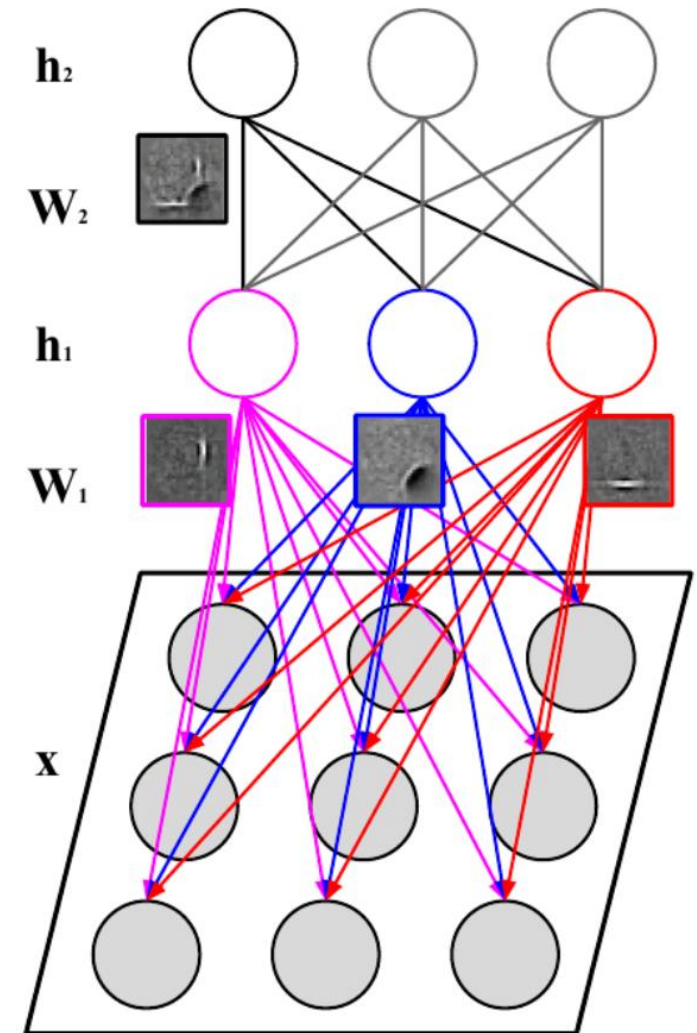o Specifically, apply Gibbs sampler for $k$ steps and approximate the gradient

$$\frac{\partial \log p(x_n | \theta)}{\partial \theta} = -\frac{\partial E(x_n, h_0 | \theta)}{\partial \theta} - \frac{\partial E(x_k, h_k | \theta)}{\partial \theta}$$



Hinton, *Training Products of Experts by Minimizing Contrastive Divergence*, Neural Computation, 2002

# Deep Belief Network

o RBMs are just one layer

o Use RBM as a building block

o Stack multiple RBMs one on top of the other
$$p(x, h_1, h_2) = p(x|h_1) \cdot p(h_1|h_2)$$

o Deep Belief Networks (DBN) are directed models
  ◦ The layers are densely connected and have a single forward flow
  ◦ This is because the RBN is directional, $p(x_i|h, \theta) = \sigma(W_{.i}x + c_i)$, namely the input argument has only variable only from below
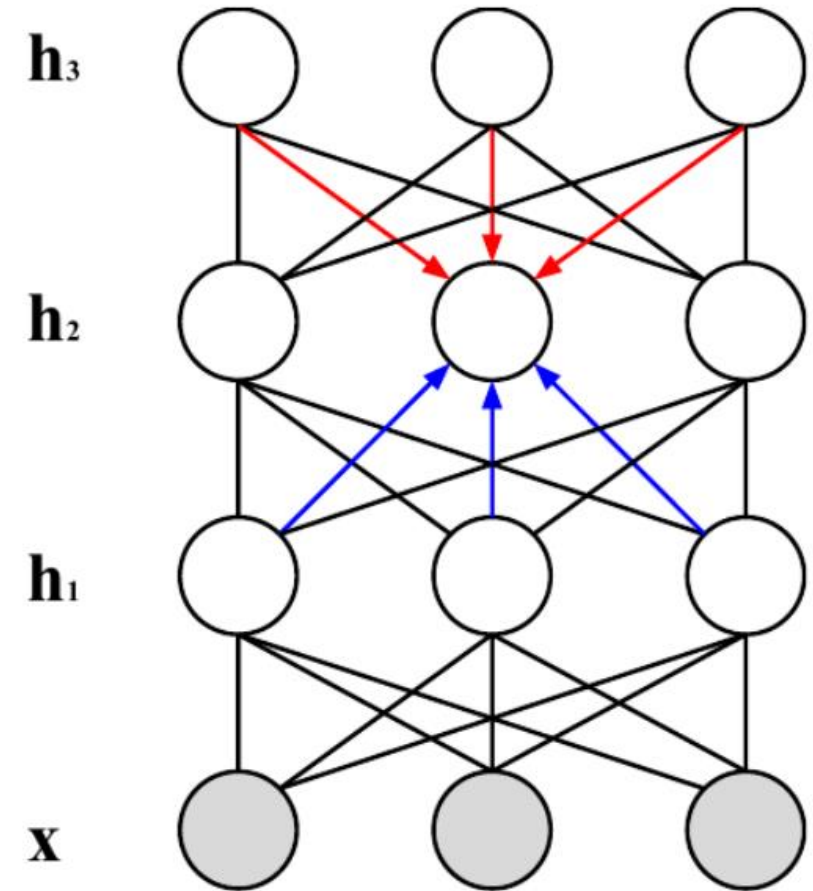
# Deep Boltzmann Machines

o Stacking layers again, but now with connection from the <span style="color:red">above</span> and from the <span style="color:blue">below</span> layers

o Since it's a Boltzmann machine, we need an energy function

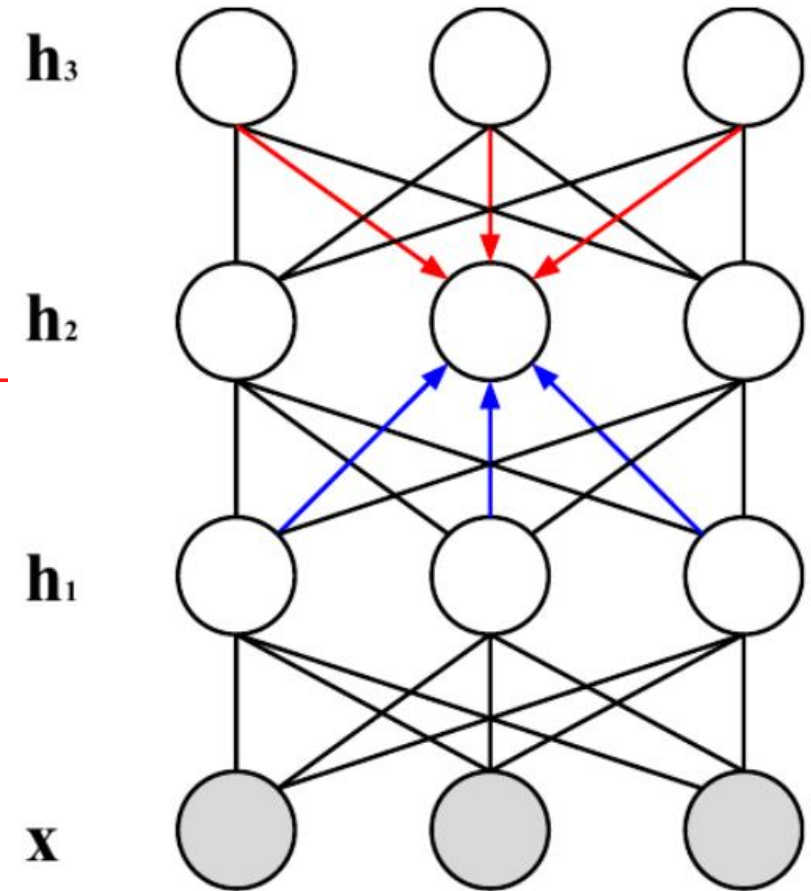$$E(x, h_1, h_2 | \theta) = x^T W_1 h_1 + h_1^T W_2 h_2 + h_2^T W_3 h_3$$

$$p(h_2^k | h_1, h_3) = \sigma(\sum_j W_1^{jk} h_1^j + \sum_l W_3^{kl} h_3^k)$$

# Deep Boltzmann Machines

o Schematically similar to Deep Belief Networks

o But, Deep Boltzmann Machines (DBM) are undirected models

◦ Belong to the Markov Random Field family

o So, two types of relationships: bottom-up and up-bottom

$$p(h_2^k|h_1, h_3) = \sigma(\sum_j W_1^{jk} h_1^j + \sum_l W_3^{kl} h_3^k)$$



$\mathbf{h_3}$

$\mathbf{h_2}$

$\mathbf{h_1}$

$\mathbf{x}$

# Training Deep Boltzmann Machines

o Computing gradients is intractable

o Instead, variational methods (mean-field) or sampling methods are used

Bayesian Modelling
Variational Inference

# Bayesian Terminology

- Observed variables $x$

- Latent variables $\theta$
  - Both unobservable model parameters $w$ and unobservable model activations $z$
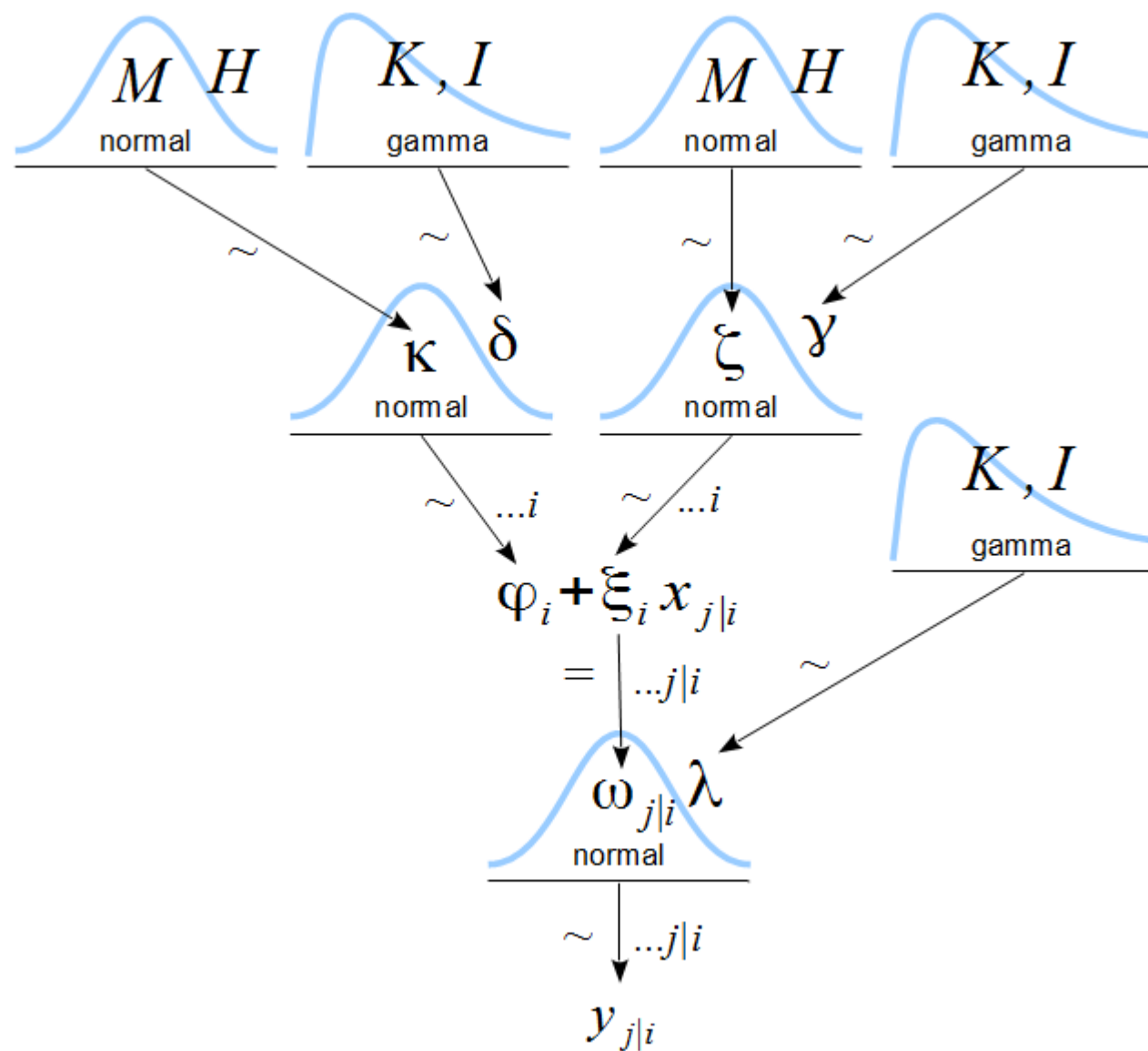  - $\theta = \{w, z\}$

- Joint probability density function (pdf): $p(x, \theta)$

- Marginal pdf: $p(x) = \int_\theta p(x, \theta)\, d\theta$

- Prior pdf ➔ marginal over input: $p(\theta) = \int_x p(x, \theta)\, dx$
  - Usually a user defined pdf

- Posterior pdf: $p(\theta|x)$

- Likelihood pdf: $p(x|\theta)$

$x$

# Bayesian Terminology

○ Posterior pdf

$$p(\theta|x) =$$  ← Conditional probability

$$= \frac{p(x,\theta)}{p(x)}$$  ← Bayes Rule

$$= \frac{p(x|\theta)\,p(\theta)}{p(x)}$$  ← Marginal probability

← $p(x\ )$ is constant

$$= \frac{p(x|\theta)\,p(\theta)}{\int_{\theta'} p(x,\theta')\,d\theta'}$$

$$\propto p(x|\theta)\,p(\theta)$$

○ Posterior Predictive pdf

$$p(y_{new}|y) = \int_{\theta} p(y_{new}|\theta)\,p(\theta|y)\,d\theta$$

# Bayesian Terminology

o Conjugate priors

◦ when posterior and prior belong to the same family, so the joint pdf is easy to compute

o Point estimate approximations of latent variables

◦ instead of computing a distribution over all possible values for the variable, compute one point only, e.g. the most likely (maximum likelihood or max a posteriori estimate)

$$\theta^* = \arg_\theta \max p(x|\theta)p(\theta)\ (MAP)$$
$$\theta^* = \arg_\theta \max p(x|\theta)\qquad (MLE)$$

◦ Quite good when the posterior distribution is peaky (low variance)

Point estimate of your neural network weight

# Bayesian Modelling

o Estimate the posterior density $p(\theta|x)$ for your training data $x$

o To do so, need to define the prior $p(\theta)$ and likelihood $p(x|\theta)$ distributions

o Once the $p(\theta|x)$ density is estimated, Bayesian Inference is possible
  ◦ $p(\theta|x)$ is a (density) function, not just a single number (point estimate)

o But how to estimate the posterior density?
  ◦ Markov Chain Monte Carlo (MCMC) ➔ Simulation-like estimation
  ◦ Variational Inference ➔ Turn estimation to optimization

# Variational Inference

○ Estimating the true posterior $p(\theta|x)$ is not always possible
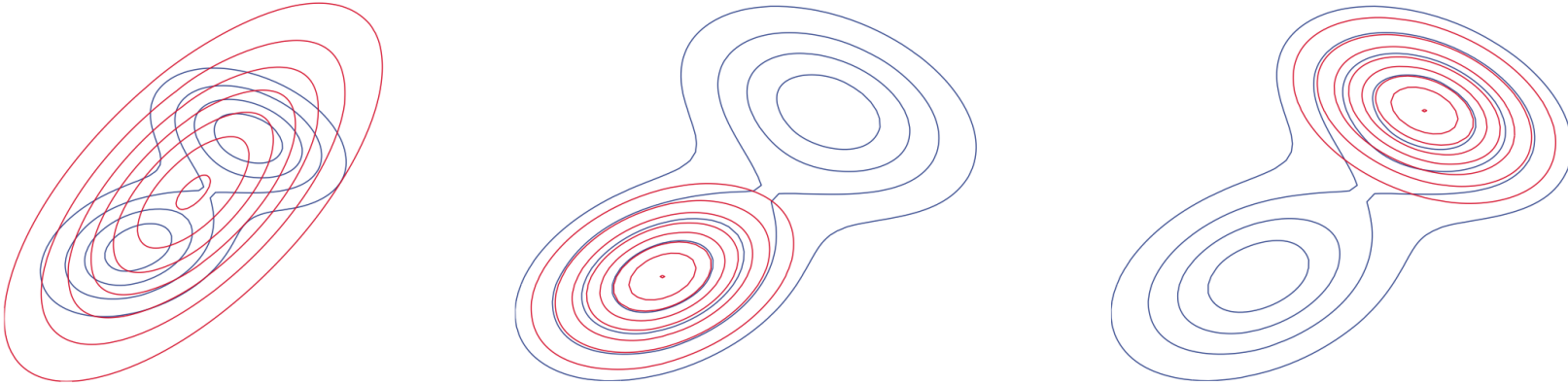  ◦ especially for complicated models like neural networks

○ Variational Inference assumes another function $q(\theta|\varphi)$ with which we want to approximate the true posterior $p(\theta|x)$
  ◦ $q(\theta|\varphi)$ is the approximate posterior
  ◦ Note that the approximate posterior does not depend on the observable variables $x$

○ We approximate by minimizing the **reverse** KL-divergence w.r.t. $\varphi$
$$\varphi^* = \arg\min_{\varphi} KL(q(\theta|\varphi)||p(\theta|x))$$

○ Turn inference into optimization

# Variational Inference (graphically)



Underestimating variance. Why?

# Variational Inference (graphically)



$p(\mathbf{z} \mid \mathbf{x})$

$q(\mathbf{z}; \boldsymbol{v})$

$\mathrm{KL}(q(\mathbf{z}; \boldsymbol{v}^*) \,\|\, p(\mathbf{z} \mid \mathbf{x}))$

$\boldsymbol{v}^*$

$\boldsymbol{v}^{\mathrm{init}}$

Underestimating variance. Why?

# Variational Inference (graphically)



$p(\mathbf{z} \mid \mathbf{x})$

$q(\mathbf{z}; \boldsymbol{v})$

$\boldsymbol{v}^*$    $\mathrm{KL}(q(\mathbf{z}; \boldsymbol{v}^*) \parallel p(\mathbf{z} \mid \mathbf{x}))$

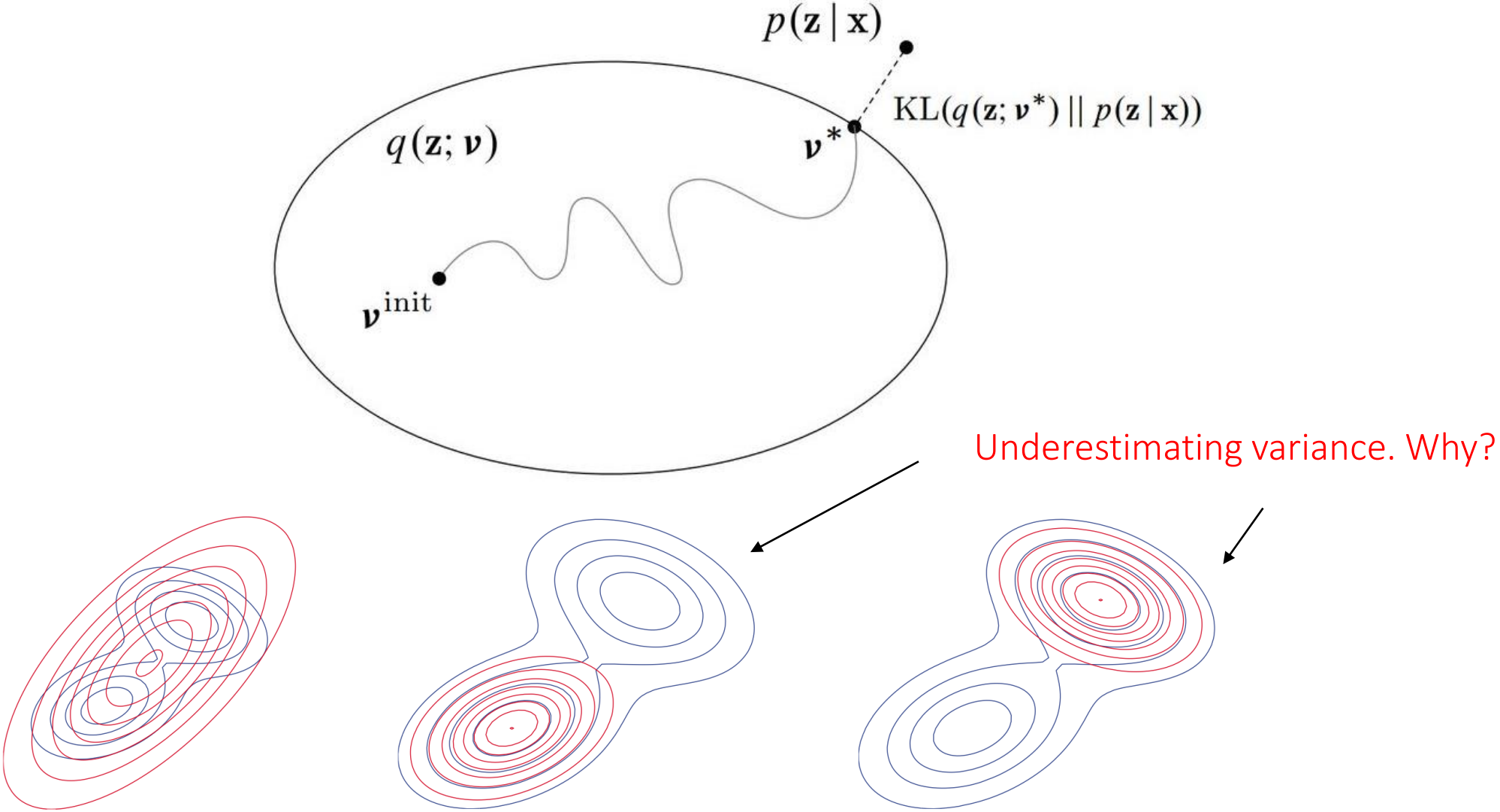$\boldsymbol{v}^{\mathrm{init}}$

How to overestimate variance?

Underestimating variance. Why?

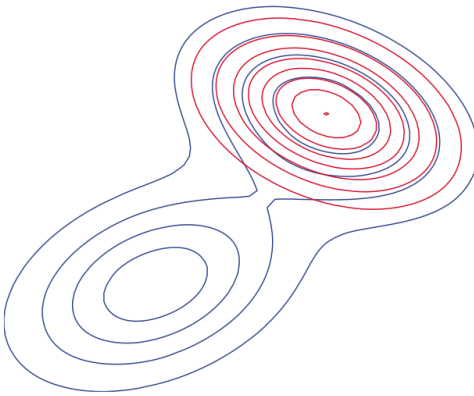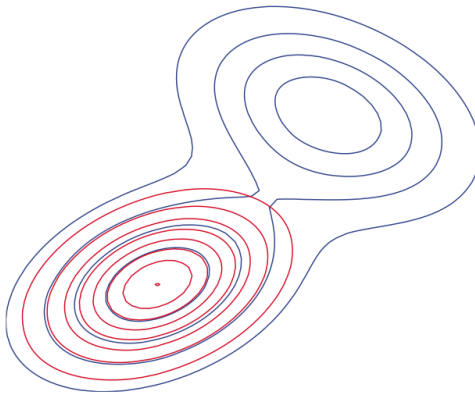# Variational Inference (graphically)



$$p(\mathbf{z} \mid \mathbf{x})$$

$$\mathrm{KL}(q(\mathbf{z}; \boldsymbol{v}^*) \parallel p(\mathbf{z} \mid \mathbf{x}))$$

$$q(\mathbf{z}; \boldsymbol{v})$$

$$\boldsymbol{v}^*$$

$$\boldsymbol{v}^{\mathrm{init}}$$

How to overestimate variance? Forward KL

Underestimating variance. Why?

# Mean-Field Approximation and CAVI Optimization

o To make the optimization of the VI easier, one can assume the latent variables are independent of each other

$$q(\theta|\varphi) = \prod_j q_j(\theta_j|\varphi_j)$$

o The optimization is often done with CAVI
- Coordinate-Ascent Variational Inference
- Initially set $\varphi$ randomly
- For each $j$ in turn you set $q_j(\theta_j|\varphi_j) = \mathbb{E}_{g_{-j}}[\log p(\theta|x)]$

# Variational Inference - Evidence Lower Bound (ELBO)

○ Given latent variables $\theta$ and the approximate posterior

$$q_\varphi(\theta) = q(\theta|\varphi)$$

○ The log marginal is

$$\log p(x) = \log \int_\theta p(x, \theta) \, d\theta$$

$$= \log \int_\theta p(x, \theta) \frac{q_\varphi(\theta)}{q_\varphi(\theta)} \, d\theta$$

$$= \log \mathbb{E}_{q_\varphi(\theta)} \left[ \frac{p(x, \theta)}{q_\varphi(\theta)} \right]$$

$$\leq \mathbb{E}_{q_\varphi(\theta)} \left[ \log \frac{p(x, \theta)}{q_\varphi(\theta)} \right]$$

$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x, \theta)] - \mathbb{E}_{q_\varphi(\theta)}[\log q_\varphi(\theta)]$$

$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x, \theta)] + H(\theta)$$

$$= \text{ELBO}_{\theta,\varphi}(x)$$

or

$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x|\theta)] - \mathbb{E}_{q_\varphi(\theta)}[\log p(\theta)]$$

$$+ \mathbb{E}_{q_\varphi(\theta)}[\log q_\varphi(\theta)]$$

$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x|\theta)] - \text{KL}(q_\varphi(\theta)||p(\theta))$$

$$= \text{ELBO}_{\theta,\varphi}(x)$$

# Variational Inference - Evidence Lower Bound (ELBO)

- Given latent variables $\theta$ and the approximate posterior

$$q_\varphi(\theta) = q(\theta|\varphi)$$

- The log marginal is

# ELBO and the marginal

o It is easy to see that the ELBO is directly related to the marginal

$$\mathrm{ELBO}_{\theta,\varphi}(x) =$$
$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x,\theta)] - \mathbb{E}_{q_\varphi(\theta)}[\log q_\varphi(\theta)]$$
$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(\theta|x)] + \mathbb{E}_{q_\varphi(\theta)}[\log p(x)] - \mathbb{E}_{q_\varphi(\theta)}[\log q_\varphi(\theta)]$$
$$= \mathbb{E}_{q_\varphi(\theta)}[\log p(x)] - KL(q_\varphi(\theta)||p(\theta|x))$$
$$= \log p(x) - KL(q_\varphi(\theta)||p(\theta|x)) \qquad \leftarrow \log p(x) \text{ does not depend on } q_\varphi(\theta)$$
$$\Rightarrow \qquad\qquad\qquad\qquad\qquad\qquad\qquad \leftarrow \mathbb{E}_{q_\varphi(\theta)}[1]=1$$
$$\log p(x) = \mathrm{ELBO}_{\theta,\varphi}(x) + KL(q_\varphi(\theta)||p(\theta|x))$$

o You can also see $\mathrm{ELBO}_{\theta,\varphi}(x)$ as Variational Free Energy

# ELBO and the marginal

○ It is easy to see that the ELBO is directly related to the marginal

$$\mathrm{ELBO}_{\theta,\varphi}(x) =$$

# ELBO interpretations

○ $\log p(x) = \mathrm{ELBO}_{\theta,\varphi}(\mathrm{x}) + KL(q_\varphi(\theta)||p(\theta|x))$

○ The log-likelihood is constant, as it does not depends on any parameter

○ Also, both $\mathrm{ELBO}_{\theta,\varphi}(\mathrm{x}) > 0$ and $KL(q_\varphi(\theta)||p(\theta|x)) > 0$

1. The higher the Variational Lower Bound $\mathrm{ELBO}_{\theta,\varphi}(\mathrm{x})$, the smaller the difference between the approximate posterior $q_\varphi(\theta)$ and the true posterior $p(\theta|x)$ → better latent representation

2. The Variational Lower Bound $\mathrm{ELBO}_{\theta,\varphi}(\mathrm{x})$ approaches the log-likelihood → better density model

# Amortized Inference

o The variational distribution $q(\theta|\varphi)$ does not depend directly on data

  ◦ Only indirectly, via minimizing its distance to the true posterior $KL(q(\theta|\varphi)||p(\theta|x))$

o So, with $q(\theta|\varphi)$ we have a major optimization problem, as the approximate posterior must approximate the whole dataset $x = [x_1, x_2, \ldots, x_N]$ jointly

o As this is obviously quite complex, one can amortize the optimization on individual data points by setting

$$q(\theta|\varphi) = q_\varphi(\theta|x)$$

o Predict model parameters $\theta$ using a $\varphi$-parameterized model of the input $x$

o Use it for parameters that depend on data, such as the latent activations

# Amortized Inference (Intuitively)

○ Originally, Variational Inference assumed that $q(\theta|\varphi)$ describes the approximate posterior of the dataset as a whole
  ◦ Think of $\theta$ not as the latent activations $z$, but only the latent model variables $w$

# Variational Autoencoders

○ Let's rewrite the ELBO a bit more explicitly

$$\text{ELBO}_{\theta, \varphi}(x) = \mathbb{E}_{q_\varphi(\theta)}[\log p(x|\theta)] - \text{KL}(q_\varphi(\theta)||\text{p}(\theta))$$
$$= \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|z)] - \text{KL}(q_\varphi(z|x)||\text{p}_\lambda(z))$$

○ Instead of $p(x|\theta)$ we have $p_\theta(x|z)$ to indicate that the model for the posterior density has weights parameterized by $\theta$ and latent model activations parameterized by $z$

○ Instead of $\text{p}(\theta)$ we have $\text{p}_\lambda(z)$, namely we put a $\lambda$-parameterized prior only on the latent activations $z$ and not the model weights

○ Instead of $q(\theta|\varphi)$ we have $q_\varphi(z|x)$ to indicate that the model approximates the posterior density of the latent activations, and the model weights are parameterized by $\varphi$
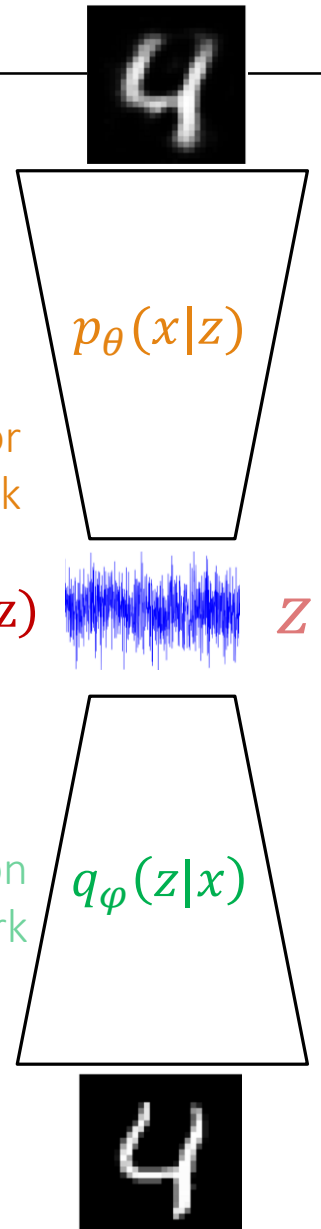
# Variational Autoencoders



- So, we have $\text{ELBO}_{\theta,\varphi}(x) = \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|z)] - \text{KL}(q_\varphi(z|x)||p_\lambda(z))$

- What if we model the densities $p_\theta(x|z)$ and $q_\varphi(z|x)$ as neural networks?

- The approximate posterior looks like a standard CovnNet (or MLP), which receives an image input $x$ and returns a feature map/latent variable $z$
  - Also known as encoder or inference network

- The likelihood term $p_\theta(x|z)$ looks like an inverted ConvNet (deconvolutions), which given a latent feature map $z$ reconstructs the input $x$
  - Also known as decoder or generator network, because it recognizes the input given the latent variable

- A difference from a standard autoencoder is we now have an opinion of what the distribution of the latents $z$ should look like, with $p_\lambda(z))$

# Training Variational Autoencoders

○ Maximize the Evidence Lower Bound (ELBO)

　◦ Or minimize the negative ELBO

$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|z)] - \mathrm{KL}(q_\varphi(z|x) \| \mathrm{p}_\lambda(\mathrm{z}))$$

○ How to we optimize the ELBO?

# Training Variational Autoencoders

o Maximize the Evidence Lower Bound (ELBO)

◦ Or minimize the negative ELBO

$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|Z)] - \mathrm{KL}(q_\varphi(Z|x)||\mathrm{p}_\lambda(Z))$$

$$= \int_z q_\varphi(z|x) \log p_\theta(x|z) \, dz - \int_z q_\varphi(z|x) \log \frac{q_\varphi(z|x)}{p_\lambda(z)} \, dz$$

o Forward propagation → compute the two terms

o The first term is an integral (expectation) that we cannot solve analytically. So, we need to sample from the pdf instead

◦ When $p_\theta(x|z)$ contains even a few nonlinearities, like in a neural network, the integral is hard to compute analytically
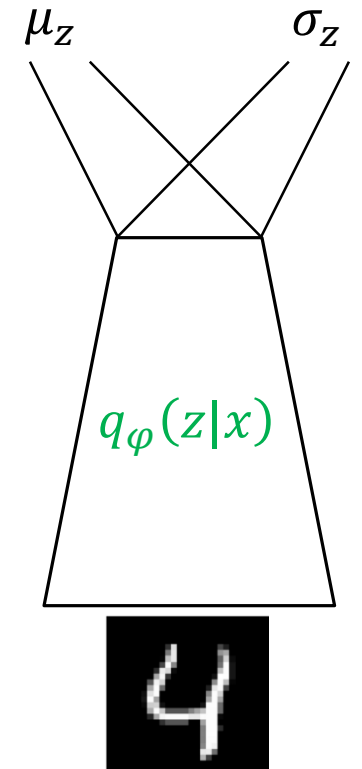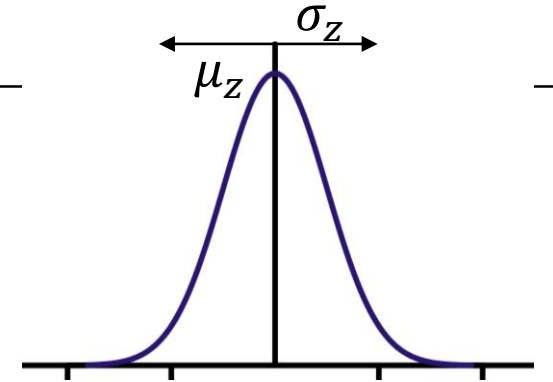
# Complex integrals

# Training Variational Autoencoders

○ Maximize the Evidence Lower Bound (ELBO)

　◦ Or minimize the negative ELBO

$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{q_\varphi(z|x)}[\log p_\theta(x|Z)] - \text{KL}(q_\varphi(Z|x)||p_\lambda(Z))$$

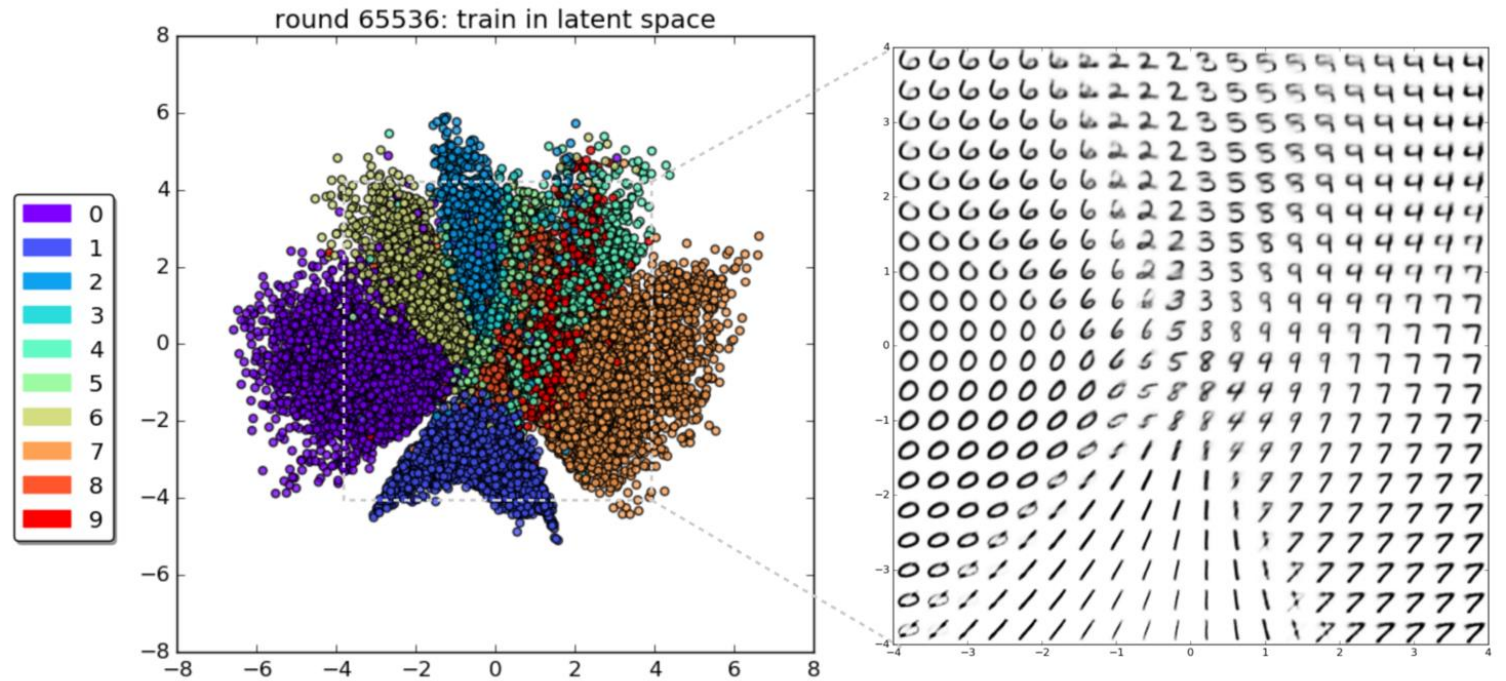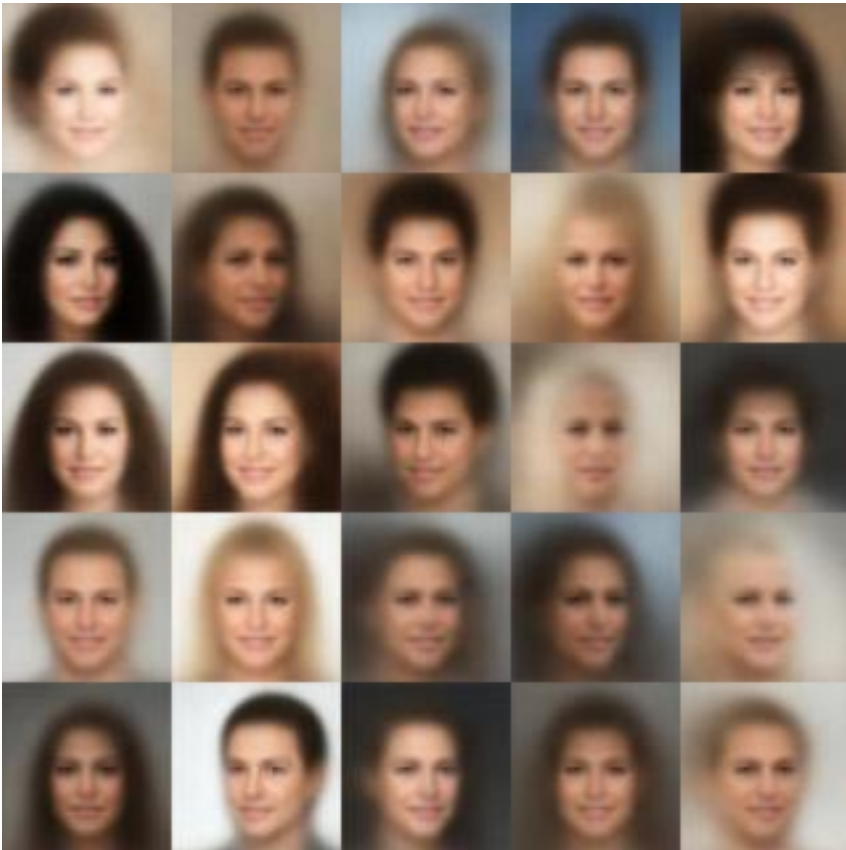$$= \int_z q_\varphi(z|x) \log p_\theta(x|z)\, dz - \int_z q_\varphi(z|x) \log \frac{q_\varphi(z|x)}{p_\lambda(z)}\, dz$$

○ Forward propagation → compute the two terms

○ The first term is an integral (expectation) that we cannot solve analytically. So, we need to sample from the pdf instead

　◦ When $p_\theta(x|z)$ contains even a few nonlinearities, like in a neural network, the integral is hard to compute analytically

○ The second term is the KL divergence between two distributions that we know

# Typical VAE



o We set the prior $p_\lambda(Z)$ to be the unit Gaussian
$$p(Z) \sim N(0, 1)$$

o We set the likelihood to be a Bernoulli for binary data
$$p(X|Z) \sim Bernoulli(\pi)$$

o We set $q_\varphi(Z|x)$ to be a neural network (MLP, ConvNet), which maps an input $x$ to the Gaussian distribution, specifically it's mean and variance

◦ $\mu_z, \sigma_z \sim q_\varphi(Z|x)$

◦ The neural network has two outputs, one is the mean $\mu_x$ and the other the $\sigma_x$, which corresponds to the covariance of the Gaussian

o We set $p_\theta(X|Z)$ to be an inverse neural network, which maps $Z$ to the Bernoulli distribution if our outputs binary (e.g. Binary MNIST)

# VAE: Interpolation in the latent space



round 65536: train in latent space

# Forward propagation in VAE

- Sample $z$ from the approximate posterior density $z \sim q_\varphi(Z|x)$

  - As $q_\varphi$ is a neural network that outputs values from a specific and known parametric pdf, e.g. a Gaussian, sampling from it is rather easy
  - Often even a single draw is enough

- Second, compute the $\log p_\theta(x|Z)$

  - As $p_\theta$ is a a neural network that outputs values from a specific and known parametric pdf, e.g. a Bernoulli for white/black pixels, computing the log-prob is easy

- Computing the ELBO is rather straightforward in the standard case

- How should we optimize the ELBO?

# Forward propagation in VAE

o Sample $z$ from the approximate posterior density $z \sim q_\varphi(Z|x)$

  ◦ As $q_\varphi$ is a neural network that outputs values from a specific and known parametric pdf, e.g. a Gaussian, sampling from it is rather easy

  ◦ Often even a single draw is enough

o Second, compute the $\log p_\theta(x|Z)$

  ◦ As $p_\theta$ is a a neural network that outputs values from a specific and known parametric pdf, e.g. a Bernoulli for white/black pixels, computing the log-prob is easy

o Computing the ELBO is rather straightforward in the standard case

o How should we optimize the ELBO? Backpropagation?

# Backward propagation in VAE

○ Backpropagation ➔ compute the gradients of
$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] - \text{KL}(q_\varphi(Z|x)||\text{p}_\lambda(Z))$$

○ $\nabla_\theta \mathcal{L} = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\nabla_\theta \log p_\theta(x|z)]$

◦ The expectation and sampling in $\mathbb{E}_{z \sim q_\varphi(Z|x)}$ does not depend on $\theta$, so no problem!

◦ Also, the $\text{KL}$ does not depend on $\theta$, so no gradient from over there!

○ $\nabla_\varphi \mathcal{L} = \nabla_\varphi \left[ \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] \right] - \nabla_\varphi \left[ \text{KL}(q_\varphi(Z|x)||\text{p}_\lambda(Z)) \right]$

# Backward propagation in VAE

o Backpropagation → compute the gradients of

$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] - \text{KL}(q_\varphi(Z|x)||\text{p}_\lambda(\text{Z}))$$

# Backward propagation in VAE

○ Backpropagation ➔ compute the gradients of
$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] - \mathrm{KL}(q_\varphi(Z|x)||\mathrm{p}_\lambda(Z))$$

○ $\nabla_\theta \mathcal{L} = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\nabla_\theta \log p_\theta(x|z)]$

◦ The expectation and sampling in $\mathbb{E}_{z \sim q_\varphi(Z|x)}$ does not depend on $\theta$, so no problem!

◦ Also, the $\mathrm{KL}$ does not depend on $\theta$, so no gradient from over there!

○ $\nabla_\varphi \mathcal{L} = \nabla_\varphi \left[ \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] \right] - \nabla_\varphi \left[ \mathrm{KL}(q_\varphi(Z|x)||\mathrm{p}_\lambda(Z)) \right]$

○ Problem?

# Backward propagation in VAE

o Backpropagation ➔ compute the gradients of

$$\mathcal{L}(\theta, \varphi) = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)] - \text{KL}(q_\varphi(Z|x)||\text{p}_\lambda(Z))$$

o $\nabla_\theta \mathcal{L} = \mathbb{E}_{z \sim q_\varphi(Z|x)}[\nabla_\theta \log p_\theta(x|z)]$

◦ The expectation and sampling in $\mathbb{E}_{z \sim q_\varphi(Z|x)}$ does not depend on $\theta$, so no problem!

◦ Also, the KL does not depend on $\theta$, so no gradient from over there!

o $\nabla_\varphi \mathcal{L} = \nabla_\varphi \left[\mathbb{E}_{z \sim q_\varphi(Z|x)}[\log p_\theta(x|z)]\right] - \nabla_\varphi\left[\text{KL}(q_\varphi(Z|x)||\text{p}_\lambda(Z))\right]$

o Problem? Sampling $z \sim q_\varphi(Z|x)$ is not differentiable ➔ no gradients

o No gradients ➔ No backprop ➔ No training! ➔ Solution?

# Solution: REINFORCE?

o So, our latent variable $Z$ is a Gaussian (in the standard VAE) represented by the mean and variance $\mu_Z, \sigma_Z$, which are the output of a neural net

o So, we can train by sampling randomly from that Gaussian
$$z \sim N(\mu_Z, \sigma_Z)$$

o Once we have that $z$, however, it's a fixed value (not a function), so we cannot backprop

o We could use, however, the REINFORCE algorithm to compute an approximation to the gradient
  ◦ High-variance gradients ➔ slow and not very effective learning

# Solution: Reparameterization trick

o Remember, we have a Gaussian output $z \sim N(\mu_Z, \sigma_Z)$

o For certain pdfs, including the Gaussian, we can rewrite their random variable $z$ as deterministic transformations of a simpler random variable $\varepsilon$

o For the Gaussian specifically, the following two formulations are equivalent

$$z \sim N(\mu_Z, \sigma_Z) \iff z = \mu_Z + \varepsilon \cdot \sigma_Z,$$

where $\varepsilon \sim N(0, 1)$ and $\mu_Z, \sigma_Z$ are deterministic values from the NN function

# Solution: Reparameterization trick

○ Instead of sampling from $z \sim N(\mu_Z, \sigma_Z)$, we sample from $\varepsilon \sim N(0, 1)$ and then we compute $z$

○ Sampling directly from $z \sim N(\mu_Z, \sigma_Z)$ leads to high-variance estimates

○ Sampling directly from $\varepsilon \sim N(0,1)$ leads to low-variance estimates
  ◦ Why low variance? Exercise for the interested reader

○ Remember: since we are sampling for $z$, we are also sampling gradients

○ More distributions beyond Gaussian possible: Laplace, Student-t, Logistic, Cauchy, Rayleight, Pareto

High-variance gradient

Low-variance gradient

# Check what is random

- Again, the latent variable is $z = \mu_Z + \varepsilon \cdot \sigma_z$

- $\mu_Z$ and $\sigma_z$ are deterministic functions (via the neural network encoder)

- $\varepsilon$ is a random variable, which comes **externally**

- The $z$ as a result is itself a random variable, because of $\varepsilon$

- However, now the randomness is not associated with the neural network and its parameters that we have to learn
  - The randomness instead comes from the external $\varepsilon$
  - The gradients flow through $\mu_Z$ and $\sigma_Z$

# Reparameterization Trick (graphically)



Original form
Reparameterised form

$\sim q(z|\phi,x)$

**Backprop**

$\partial f/\partial z_j$    $z$    $= g(\phi,x,\varepsilon)$

$\partial f/\partial \varphi_i$

$\simeq \partial L/\partial \varphi_i$

$\varepsilon \sim p(\varepsilon)$

◇ : Deterministic node

● : Random node

[Kingma, 2013]
[Bengio, 2013]
[Kingma and Welling 2014]
[Rezende et al 2014]

# VAE Training Pseudocode

**Data:**
$\mathcal{D}$: Dataset
$q_\phi(\mathbf{z}|\mathbf{x})$: Inference model
$p_\theta(\mathbf{x}, \mathbf{z})$: Generative model
**Result:**
$\theta, \phi$: Learned parameters

$(\theta, \phi) \leftarrow$ Initialize parameters
**while** *SGD not converged* **do**
$\quad \mathcal{M} \sim \mathcal{D}$ (Random minibatch of data)
$\quad \epsilon \sim p(\epsilon)$ (Random noise for every datapoint in $\mathcal{M}$)
$\quad$ Compute $\tilde{\mathcal{L}}_{\theta,\phi}(\mathcal{M}, \epsilon)$ and its gradients $\nabla_{\theta,\phi}\tilde{\mathcal{L}}_{\theta,\phi}(\mathcal{M}, \epsilon)$
$\quad$ Update $\theta$ and $\phi$ using SGD optimizer
**end**

The ELBO's gradients

# VAE for NLP

> " i want to talk to you . "
> "i want to be with you . "
> "i do n't want to be with you . "
> i do n't want to be with you .
> she did n't want to be with him .

> he was silent for a long moment .
> he was silent for a moment .
> it was quiet for a moment .
> it was dark and cold .
> there was a pause .
> it was my turn .

Figure 2.D.2: An application of VAEs to interpolation between pairs of sentences, from [Bowman et al., 2015]. The intermediate sentences are grammatically correct, and the topic and syntactic structure are typically locally consistent.

# VAE for Image Resynthesis



*Smile vector:*
mean smiling faces –
mean no-smile faces

**Latent space arithmetic**

Figure 2.D.3: VAEs can be used for image re-synthesis. In this example by White [2016], an original image (left) is modified in a latent space in the direction of a *smile vector*, producing a range of versions of the original, from smiling to sadness. Notice how changing the image along a single vector in latent space, modifies the image in many subtle and less-subtle ways in pixel space.

# VAE for designing chemical compounds



Figure 2.D.1: Example application of a VAE in [Gómez-Bombarelli et al., 2016]: design of new molecules with desired chemical properties. (a) A latent continuous representation $z$ of molecules is learned on a large dataset of molecules. (b) This continuous representation enables gradient-based search of new molecules that maximizes some chosen desired chemical property given by objective function $f(z)$.
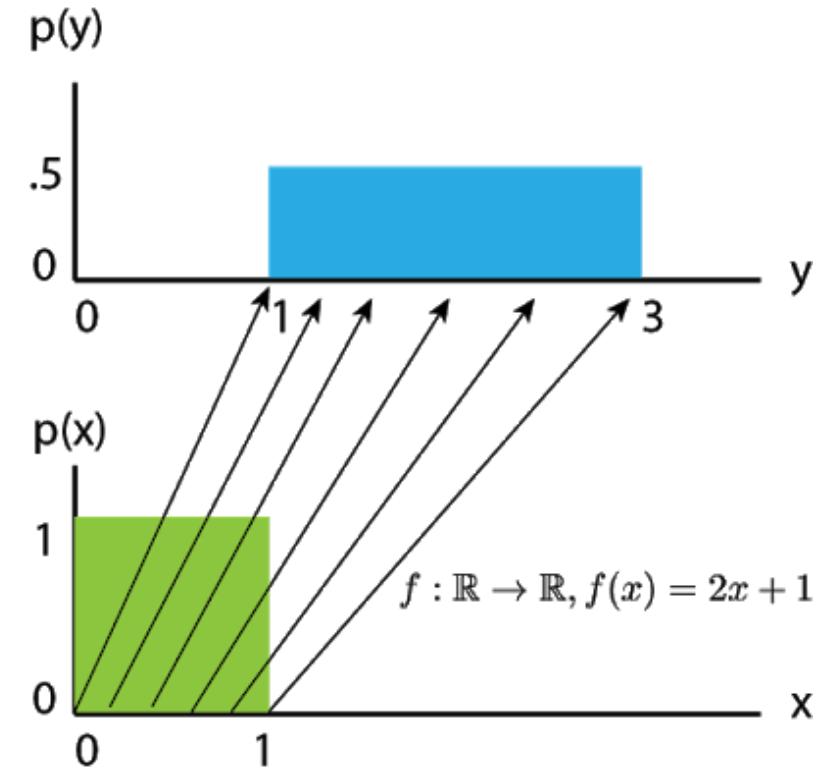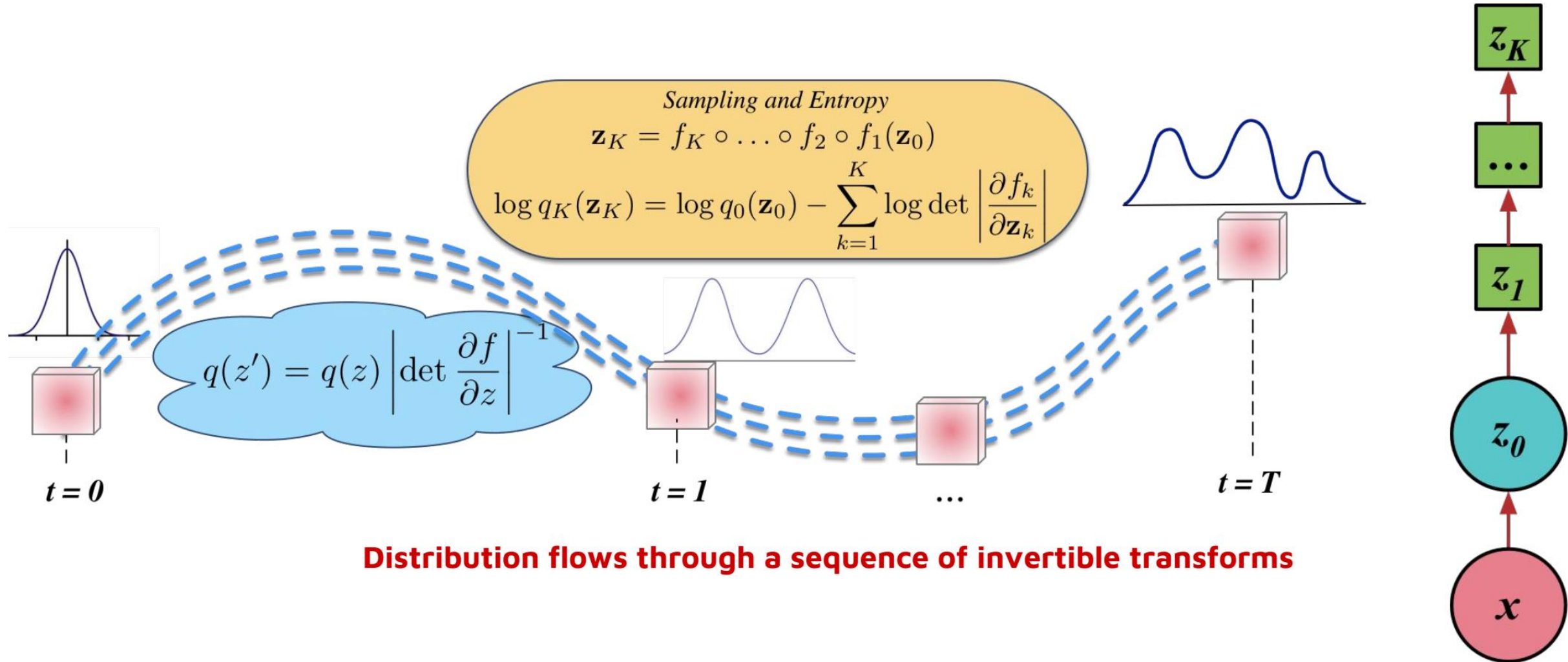
# Normalizing Flows

o Using simple pdfs, like a Gaussian, for the approximate posterior limits the expressivity of the model

o Better make sure the approximate posterior comes from a class of models that can <u>even</u> contain the true posterior

o Use a series of $K$ invertible transformations to construct the approximate posterior
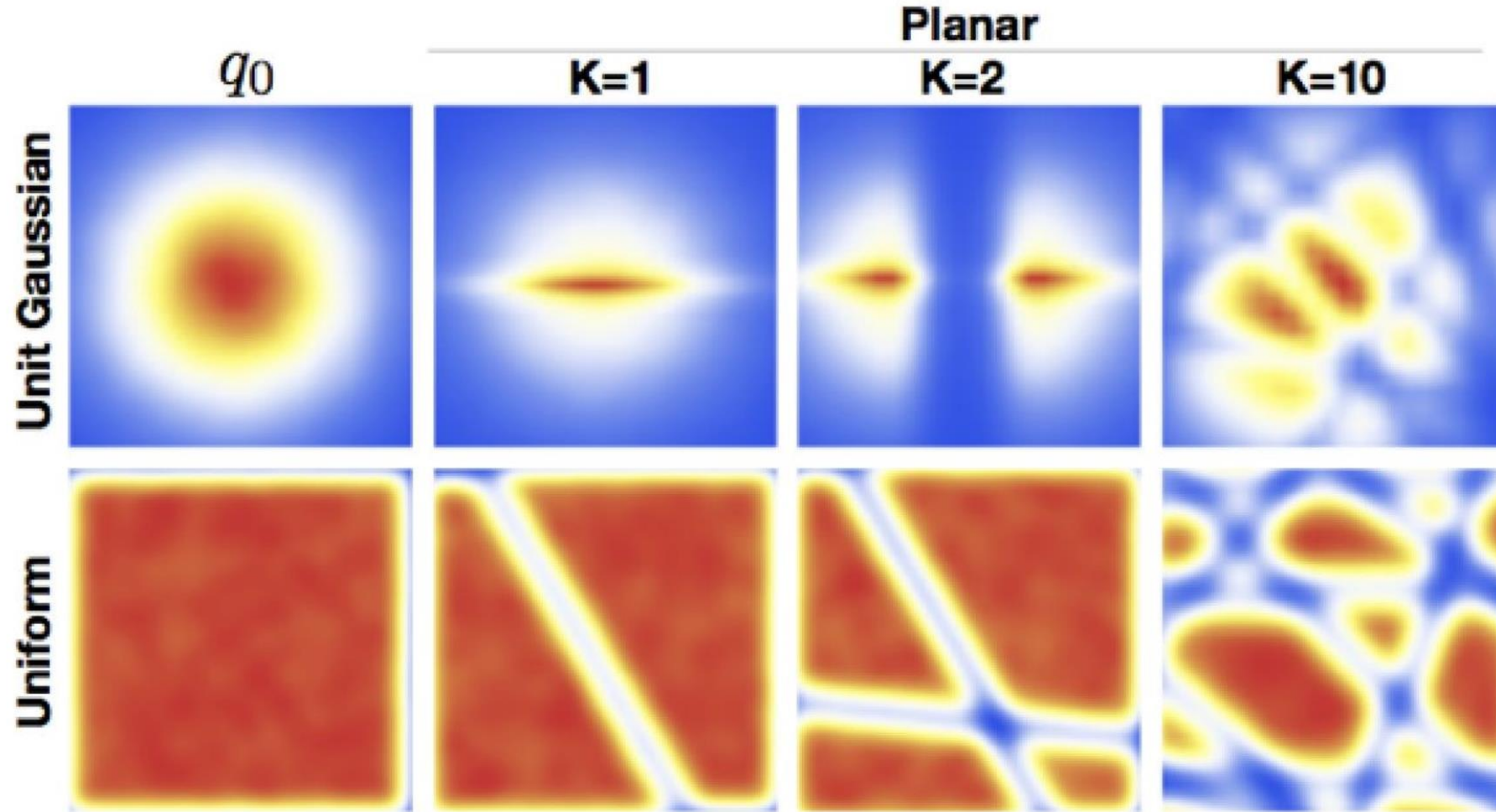  ◦ $z_k = f_k \circ f_{k-1} \circ \cdots f_1(z_0)$
  ◦ Rule of change for variables



Changing from the $x$ variable to $y$ using the transformation $y = f(x) = 2x + 1$

# Normalizing Flows

**Sampling and Entropy**

$$\mathbf{z}_K = f_K \circ \ldots \circ f_2 \circ f_1(\mathbf{z}_0)$$

$$\log q_K(\mathbf{z}_K) = \log q_0(\mathbf{z}_0) - \sum_{k=1}^{K} \log \det \left| \frac{\partial f_k}{\partial \mathbf{z}_k} \right|$$

$$q(z') = q(z) \left| \det \frac{\partial f}{\partial z} \right|^{-1}$$

$t = 0$        $t = 1$        $\ldots$        $t = T$

**Distribution flows through a sequence of invertible transforms**

# Normalizing Flows

https://www.shakirm.com/slides/DeepGenModelsTutorial.pdf

# Normalizing Flows on Non-Euclidean Manifolds



Figure 1: Left: Construction of a complex density on $\mathbf{S}^n$ by first projecting the manifold to $\mathbf{R}^n$, transforming the density and projecting it back to $\mathbf{S}^n$. Right: Illustration of transformed ($\mathbf{S}^2 \to \mathbf{R}^2$) densities corresponding to an uniform density on the sphere. Blue: empirical density (obtained by Monte Carlo); Red: Analytical density from equation (4); Green: Density computed ignoring the intrinsic dimensionality of $\mathbf{S}^n$.
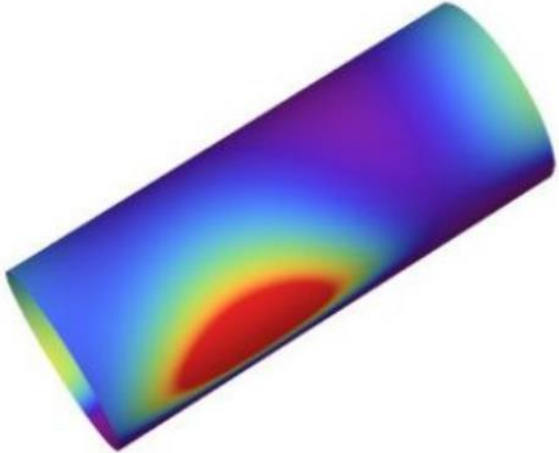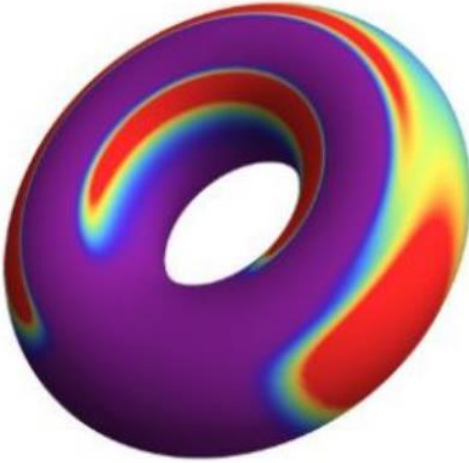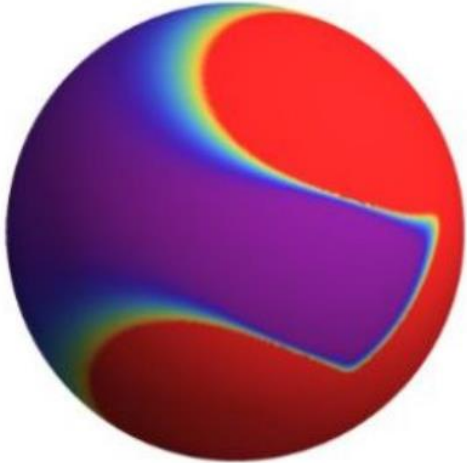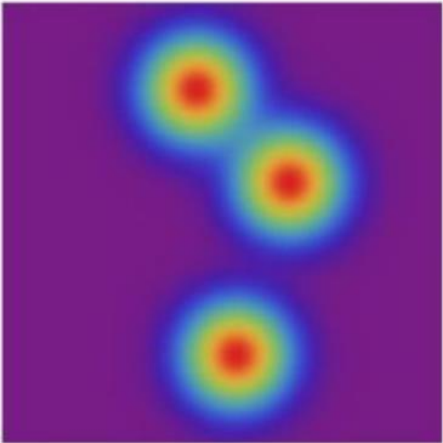
$$\log q_K(\mathbf{z}_K) = \log q_0(\mathbf{z}_0) - \frac{1}{2} \sum_{k=1}^{K} \log \det \left| \mathbf{J}_\phi^\top \mathbf{J}_\phi \right|$$

*Gemici et al., 2016*

https://www.shakirm.com/slides/DeepGenModelsTutorial.pdf

# Normalizing Flows on Non-Euclidean Manifolds

# Summary

- Gentle intro to Bayesian Modelling and Variational Inference
- Restricted Boltzmann Machines
- Deep Boltzmann Machines
- Deep Belief Network
- Contrastive Divergence
- Variational Autoencoders
- Normalizing Flows