

# Normalising Flows

Philip Schulz and Wilker Aziz

[https:](https://github.com/philschulz/VITutorial)

[//github.com/philschulz/VITutorial](https://github.com/philschulz/VITutorial)

## The problem with Standard Distributions

Normalising Flows

Use Case 1: Density Estimation

Use Case 2: Inference (sampling)

Summary

## The problem with Standard Distributions

Normalising Flows

Use Case 1: Density Estimation

Use Case 2: Inference (sampling)

Summary

# The Case of Pictures

Have you modeled pixels as Gaussian variables? Do we really believe that the pixels follow a Gaussian distribution?





# The case of Word Embeddings

# The case of Word Embeddings

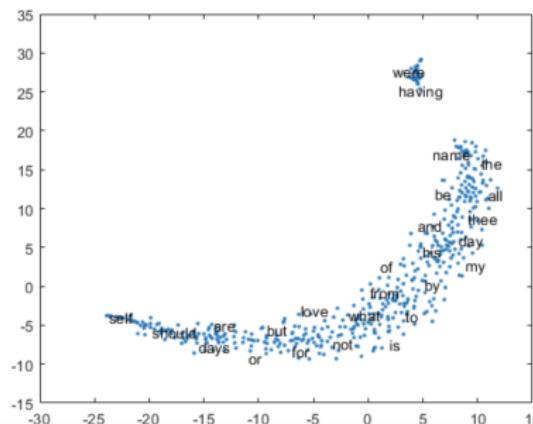


Figure: <https://it.mathworks.com/help/textanalytics/ref/trainwordembedding.html>

# Posterior Approximations

We often use exponential families to approximate posteriors. Thus we assume unimodal posteriors. Is that realistic?

# Posterior Approximations

We often use exponential families to approximate posteriors. Thus we assume unimodal posteriors. Is that realistic?

## Counter example

Gaussian mixture model

The problem with Standard Distributions

## Normalising Flows

Use Case 1: Density Estimation

Use Case 2: Inference (sampling)

Summary

## Recap: Reparametrisation

Express the density of a variable  $Y$  in terms of the density of a variable  $X$ . Assume that a differentiable, invertible mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  exists.

$$h(x) = y$$

## Recap: Reparametrisation

Express the density of a variable  $Y$  in terms of the density of a variable  $X$ . Assume that a differentiable, invertible mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  exists.

$$h(x) = y$$

$$p(y) = p(h^{-1}(y)) |\det J_{h^{-1}}(y)| = p(x) |\det J_{h^{-1}}(y)|$$

## Recap: Reparametrisation

Express the density of a variable  $Y$  in terms of the density of a variable  $X$ . Assume that a differentiable, invertible mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  exists.

$$h(x) = y$$

$$p(y) = p(h^{-1}(y)) |\det J_{h^{-1}}(y)| = p(x) |\det J_{h^{-1}}(y)|$$

$$p(x) = p(h(x)) |\det J_h(x)| = p(y) |\det J_h(x)|$$

## Recap: Reparametrisation

Express the density of a variable  $Y$  in terms of the density of a variable  $X$ . Assume that a differentiable, invertible mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  exists.

$$h(x) = y$$

$$p(y) = p(h^{-1}(y))|\det J_{h^{-1}}(y)| = p(x)|\det J_{h^{-1}}(y)|$$

$$p(x) = p(h(x))|\det J_h(x)| = p(y)|\det J_h(x)|$$

## The Challenge

The mapping  $h$  (or its inverse) needs to be defined.

# Normalising Flows

## Approach

Let's learn the transformation  $h$  (or its inverse).

# Normalising Flows

## Approach

Let's learn the transformation  $h$  (or its inverse).

## Problem

If we want  $p(y)$ , we need to provide  $|\det J_{h^{-1}}(y)|$  **in the forward pass**. But that's hard!

We are going to devise ways to get  $|\det J_{h^{-1}}(y)|$ .

# Normalising Flows

## Core Idea

Decompose mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  into

$$h = h_1 \circ h_2 \circ \dots \circ h_K .$$

# Normalising Flows

## Core Idea

Decompose mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  into

$$h = h_1 \circ h_2 \circ \dots \circ h_K .$$

Now we can learn  $K$  mappings with simple jacobian determinants.

# Normalising Flows

## Core Idea

Decompose mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  into

$$h = h_1 \circ h_2 \circ \dots \circ h_K .$$

Now we can learn  $K$  mappings with simple jacobian determinants.

$$h^{-1} = h_K^{-1} \circ h_{K-1}^{-1} \circ \dots \circ h_1^{-1}$$

# Normalising Flows

## Core Idea

Decompose mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  into

$$h = h_1 \circ h_2 \circ \dots \circ h_K .$$

Now we can learn  $K$  mappings with simple jacobian determinants.

$$h^{-1} = h_K^{-1} \circ h_{K-1}^{-1} \circ \dots \circ h_1^{-1}$$

$$p(x) = p(y) |\det J_{h_1}(y^{(1)})| |\det J_{h_2}(y^{(2)})| \dots |\det J_{h_K}(x)|$$

# Normalising Flows

## Core Idea

Decompose mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  into

$$h = h_1 \circ h_2 \circ \dots \circ h_K .$$

Now we can learn  $K$  mappings with simple jacobian determinants.

$$h^{-1} = h_K^{-1} \circ h_{K-1}^{-1} \circ \dots \circ h_1^{-1}$$

$$p(x) = p(y) |\det J_{h_1}(y^{(1)})| |\det J_{h_2}(y^{(2)})| \dots |\det J_{h_K}(x)|$$

$$p(y) = p(x) \left| \det J_{h_1^{-1}}(y^{(K-1)}) \right| \left| \det J_{h_2^{-1}}(y^{(K-2)}) \right| \dots \left| \det J_{h_1^{-1}}(y) \right|$$

The problem with Standard Distributions

Normalising Flows

Use Case 1: Density Estimation

Use Case 2: Inference (sampling)

Summary

# Normalising Flows: Density Estimation

## Setting

Our data  $x$  has unknown continuous density  $p(x)$ .  
We can therefore not handcraft a likelihood.

# Normalising Flows: Density Estimation

## Setting

Our data  $x$  has unknown continuous density  $p(x)$ .  
We can therefore not handcraft a likelihood.

## Goal

Transform known variable  $x$  into  $\epsilon = h(x)$  and  
express the likelihood as

# Normalising Flows: Density Estimation

## Setting

Our data  $x$  has unknown continuous density  $p(x)$ .  
We can therefore not handcraft a likelihood.

## Goal

Transform known variable  $x$  into  $\epsilon = h(x)$  and  
express the likelihood as

$$p(x) = p(\epsilon) |\det J_h(x)|$$

# Normalising Flows: Density Estimation

## Setting

Our data  $x$  has unknown continuous density  $p(x)$ .  
We can therefore not handcraft a likelihood.

## Goal

Transform known variable  $x$  into  $\epsilon = h(x)$  and  
express the likelihood as

$$\begin{aligned} p(x) &= p(\epsilon) |\det J_h(x)| \\ &= p(\epsilon) |\det J_{h_1}(\epsilon^{(1)})| |\det J_{h_2}(\epsilon^{(2)})| \dots |\det J_{h_K}(x)| \end{aligned}$$

# Normalising Flows: Density Estimation

## Setting

Our data  $x$  has unknown continuous density  $p(x)$ .  
We can therefore not handcraft a likelihood.

## Goal

Transform known variable  $x$  into  $\epsilon = h(x)$  and  
express the likelihood as

$$\begin{aligned} p(x) &= p(\epsilon) |\det J_h(x)| \\ &= p(\epsilon) |\det J_{h_1}(\epsilon^{(1)})| |\det J_{h_2}(\epsilon^{(2)})| \dots |\det J_{h_K}(x)| \\ &= p(h_1(\epsilon^{(1)})) |\det J_{h_1}(\epsilon^{(1)})| \dots |\det J_{h_K}(x)| \end{aligned}$$

# 2-step Flow

$$p(x) = p(\epsilon) \left| \det J_{h_1^{-1}} (\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}} (x) \right|$$

## 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}} \left( \epsilon^{(1)} \right) \right| \left| \det J_{h_2^{-1}} (x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}} \left( h_2^{-1}(x) \right) \right| \left| \det J_{h_2^{-1}} (x) \right| \end{aligned}$$

## 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

The transformations  $h_1^{-1}$  and  $h_2^{-1}$  are learned by backprop. The determinants need to be computed analytically.

# Designing a Transformation

Assume:  $x_i = (x_{i1}, x_{i2}, \dots, x_{iM})$ . Then factorise the density according to the chain rule.

$$\log p(x_i|\theta) = \sum_{j=1}^M \log p(x_{ij}|x_{i,<j}\theta)$$

Next assume an invertible mapping  $h(x_{ij}) = \epsilon_{ij}$ .

## Simple Mapping

$$\begin{aligned} h(x) &= \epsilon \\ h^{-1}(\epsilon) &= x \end{aligned}$$

# Designing a Transformation

Assume:  $x_i = (x_{i1}, x_{i2}, \dots, x_{iM})$ . Then factorise the density according to the chain rule.

$$\log p(x_i|\theta) = \sum_{j=1}^M \log p(x_{ij}|x_{i,<j}\theta)$$

Next assume a mapping  $h(x_{ij}) = \epsilon_{ij}$ .

## Flow Mapping

$$h_1 \circ h_2 \circ \dots \circ h_K(x) = \epsilon$$

$$h_K^{-1} \circ h_{K-1}^{-1} \circ \dots \circ h_1^{-1}(\epsilon) = x$$

# Designing a Transformation

MADE (Germain et al., 2015)

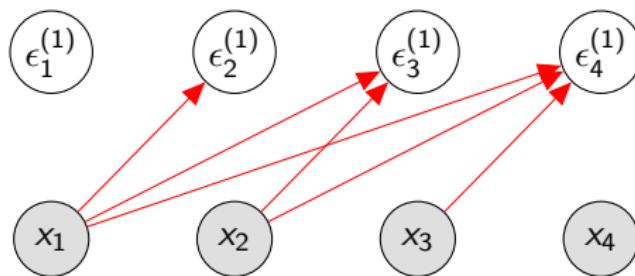
An autoregressive network that takes constant time.  
Its connectivity matrix is lower-triangular.

$$\begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 \\ 1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 1 & 1 & \cdots & 1 & 0 \end{bmatrix}$$

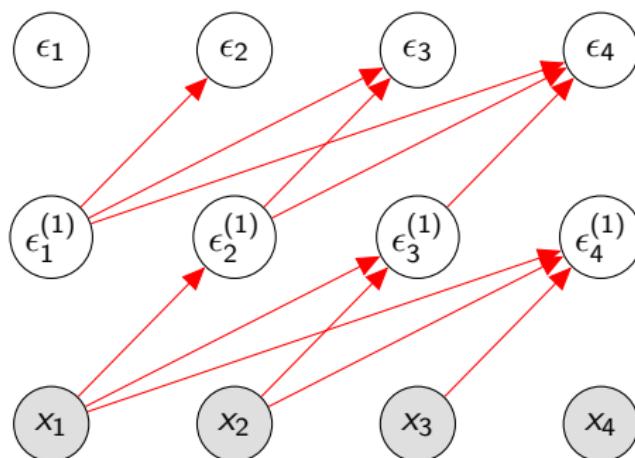
# Designing a Transformation



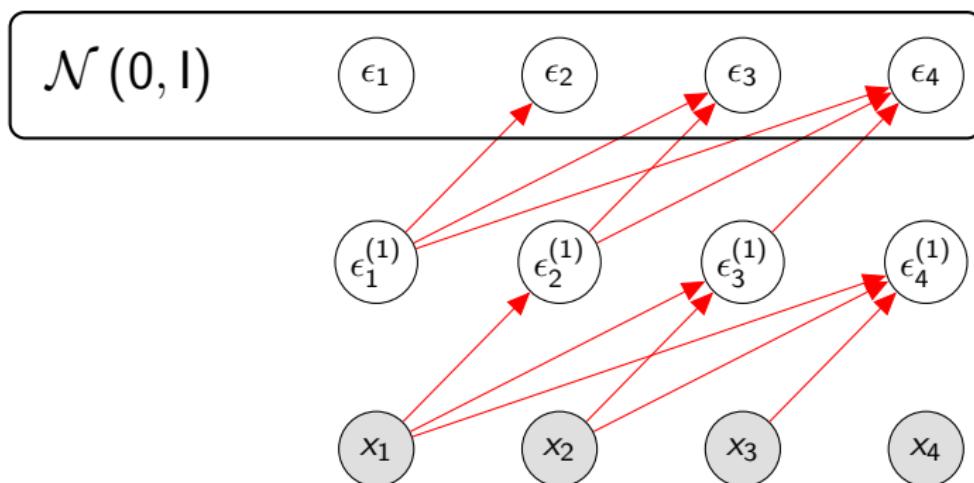
# Designing a Transformation



# Designing a Transformation



# Designing a Transformation



# Designing a Transformation

We use a MADE  $g_\theta^{(2)}$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = g_\theta^{(2)}(x_{<j})$ . Then we apply the first transformation.

$$\epsilon_j^{(1)} = h_2^{-1}(x)_j = \frac{x_j - \mu(x_{<j})}{\sigma(x_{<j})}$$

# Designing a Transformation

We use a MADE  $g_\theta^{(2)}$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = g_\theta^{(2)}(x_{<j})$ . Then we apply the first transformation.

$$\begin{aligned}\epsilon_j^{(1)} &= h_2^{-1}(x)_j = \frac{x_j - \mu(x_{<j})}{\sigma(x_{<j})} \\ \epsilon^{(1)} &= h_2^{-1}(x) = \frac{x - \mu}{\sigma}\end{aligned}$$

# Designing a Transformation

We use a MADE  $g_\theta^{(2)}$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = g_\theta^{(2)}(x_{<j})$ . Then we apply the first transformation.

$$\begin{aligned}\epsilon_j^{(1)} &= h_2^{-1}(x)_j = \frac{x_j - \mu(x_{<j})}{\sigma(x_{<j})} \\ \epsilon^{(1)} &= h_2^{-1}(x) = \frac{x - \mu}{\sigma}\end{aligned}$$

The Jacobian is

$$J_{h_2^{-1}}(x) =$$

# Designing a Transformation

We use a MADE  $g_\theta^{(2)}$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = g_\theta^{(2)}(x_{<j})$ . Then we apply the first transformation.

$$\begin{aligned}\epsilon_j^{(1)} &= h_2^{-1}(x)_j = \frac{x_j - \mu(x_{<j})}{\sigma(x_{<j})} \\ \epsilon^{(1)} &= h_2^{-1}(x) = \frac{x - \mu}{\sigma}\end{aligned}$$

The Jacobian is

$$J_{h_2^{-1}}(x) = |\sigma^{-1} + J_{\frac{-\mu}{\sigma}}(x)|$$

# Designing a Transformation

Define  $\alpha_{lj} = \frac{d}{dx_l} \frac{-\mu_j}{\sigma_j}$ .

$$J_{h_K^{-1}}(x) = |\sigma|^{-1} + J_{\frac{-\mu}{\sigma}}(x) =$$

# Designing a Transformation

Define  $\alpha_{lj} = \frac{d}{dx_l} \frac{-\mu_j}{\sigma_j}$ .

$$J_{h_K^{-1}}(x) = |\sigma^{-1}| + J_{\frac{-\mu}{\sigma}}(x) =$$

$$\begin{bmatrix} \sigma_{11}^{-1} & 0 & \cdots & 0 & 0 \\ 0 & \sigma_{22}^{-1} & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & \sigma_{mm}^{-1} \end{bmatrix}$$

# Designing a Transformation

Define  $\alpha_{lj} = \frac{d}{dx_l} \frac{-\mu_j}{\sigma_j}$ .

$$J_{h_K^{-1}}(x) = |\sigma|^{-1} + J_{\frac{-\mu}{\sigma}}(x) =$$

$$\begin{bmatrix} \sigma_{11}^{-1} & 0 & \cdots & 0 & 0 \\ 0 & \sigma_{22}^{-1} & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & \sigma_{mm}^{-1} \end{bmatrix} + \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 \\ \alpha_{21} & 0 & \cdots & 0 & 0 \\ \alpha_{31} & \alpha_{32} & \cdots & 0 & 0 \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \alpha_{m1} & \alpha_{m2} & \cdots & \alpha_{m,m-1} & 0 \end{bmatrix}$$

# Designing a Transformation

## Simple Jacobian Determinant

$$\left| \det J_{h_2^{-1}}(x) \right| = \prod_{j=1}^M \sigma_j^{-1}$$

# Designing a Transformation

## Simple Jacobian Determinant

$$\left| \det J_{h_2^{-1}}(x) \right| = \prod_{j=1}^M \sigma_j^{-1}$$

In practice we work with the log-likelihood.

$$\log \left| \det J_{h_2^{-1}}(x) \right| = - \sum_{j=1}^M \log \sigma_j$$

# 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

## 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

$$\log p(x) = \log p(h_1^{-1}(h_2^{-1}(x)))$$

## 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

$$\log p(x) = \log p(h_1^{-1}(h_2^{-1}(x))) - \sum_{j=1}^M \log \sigma_j^{(2)} - \sum_{j=1}^M \log \sigma_j^{(1)}$$

## 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

$$\log p(x) = \log p(h_1^{-1}(h_2^{-1}(x))) - \sum_{j=1}^M \log \sigma_j^{(2)} - \sum_{j=1}^M \log \sigma_j^{(1)}$$

$$\epsilon^{(1)} = h_2^{-1} = \frac{x - \mu^{(2)}}{\sigma^{(2)}} \text{ where } [\mu^{(2)}, \sigma^{(2)}] = g^{(2)}(x)$$

# 2-step Flow

$$\begin{aligned} p(x) &= p(\epsilon) \left| \det J_{h_1^{-1}}(\epsilon^{(1)}) \right| \left| \det J_{h_2^{-1}}(x) \right| \\ &= p(h_1^{-1}(h_2^{-1}(x))) \left| \det J_{h_1^{-1}}(h_2^{-1}(x)) \right| \left| \det J_{h_2^{-1}}(x) \right| \end{aligned}$$

$$\log p(x) = \log p(h_1^{-1}(h_2^{-1}(x))) - \sum_{j=1}^M \log \sigma_j^{(2)} - \sum_{j=1}^M \log \sigma_j^{(1)}$$

$$\epsilon^{(1)} = h_2^{-1} = \frac{x - \mu^{(2)}}{\sigma^{(2)}} \text{ where } [\mu^{(2)}, \sigma^{(2)}] = g^{(2)}(x)$$

$$\epsilon = h_1^{-1} = \frac{\epsilon^{(1)} - \mu^{(1)}}{\sigma^{(1)}} \text{ where } [\mu^{(1)}, \sigma^{(1)}] = g^{(1)}(\epsilon^{(1)})$$

# Intermediate Summary

- ▶ NFs map transform complex distributions to simpler ones (or vice versa)
- ▶ Use in density estimation for complex distributions
- ▶ Jacobian needs to be carefully designed
- ▶ Sampling is slow because sequential

## The problem with Standard Distributions

Normalising Flows

Use Case 1: Density Estimation

Use Case 2: Inference (sampling)

Summary

# Setting

We have a generative model  $p(x|z)$ . We want to approximate the posterior  $p(z|x)$  using an amortized variational distribution  $q(z|x)$  computed by a neural net.

# Setting

We have a generative model  $p(x|z)$ . We want to approximate the posterior  $p(z|x)$  using an amortized variational distribution  $q(z|x)$  computed by a neural net.

## Goal

We want a complex, multimodal approximate posterior  $q(z|x)$ .

# Normalising Flows: Inference

$$\begin{aligned}\text{ELBO} &= -\text{KL}(p(z|x) \parallel q(z|x)) \\ &= \mathbb{E}_{q(z|\lambda)} [\log p(x|z)] - \text{KL}(q(z|\lambda) \parallel p(z)) \\ &= \underbrace{\mathbb{E}_{q(\epsilon)} [\log p(x|h^{-1}(\epsilon))]}_{\text{sample } z} - \underbrace{\text{KL}(q(z|\lambda) \parallel p(z))}_{\text{assess density}}\end{aligned}$$

## Simple Mapping

$$\begin{aligned}h(z) &= \epsilon \text{ s.t. } \epsilon \perp \lambda \\ h^{-1}(\epsilon) &= z\end{aligned}$$

# Normalising Flows: Inference

$$\begin{aligned} -\text{KL}(q(z|x) \parallel p(z|x)) &\propto \text{ELBO} = \\ &= \mathbb{E}_{q(z|\lambda)} [\log p(x|z)] - \text{KL}(q(z|\lambda) \parallel p(z)) \\ &= \underbrace{\mathbb{E}_{q(\epsilon)} [\log p(x|h^{-1}(\epsilon))]}_{\text{sample } z} - \underbrace{\text{KL}(q(z|\lambda) \parallel p(z))}_{\text{assess density}} \end{aligned}$$

## Flow Mapping

$$\begin{aligned} h_1 \circ h_2 \circ \dots \circ h_K(z) &= \epsilon \text{ s.t. } \epsilon \perp \lambda \\ h_K^{-1} \circ h_{K-1}^{-1} \circ \dots \circ h_1^{-1}(\epsilon) &= z \end{aligned}$$

# 2-step Flow

$$q(z^{(2)}) = q(\epsilon) \left| \det J_{h_1}(z^{(1)}) \right| \left| \det J_{h_2}(z^{(2)}) \right|$$

# 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) |\det J_{h_1}(z^{(1)})| |\det J_{h_2}(z^{(2)})| \\ &= q(h_1(h_2(z^{(2)}))) |\det J_{h_1}(h_2(z^{(2)}))| |\det J_{h_2}(z^{(2)})| \end{aligned}$$

## 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) |\det J_{h_1}(z^{(1)})| |\det J_{h_2}(z^{(2)})| \\ &= q(h_1(h_2(z^{(2)}))) |\det J_{h_1}(h_2(z^{(2)}))| |\det J_{h_2}(z^{(2)})| \end{aligned}$$

The transformations  $h_1^{-1}$  and  $h_2^{-1}$  are learned by backprop. The determinants need to be computed analytically.

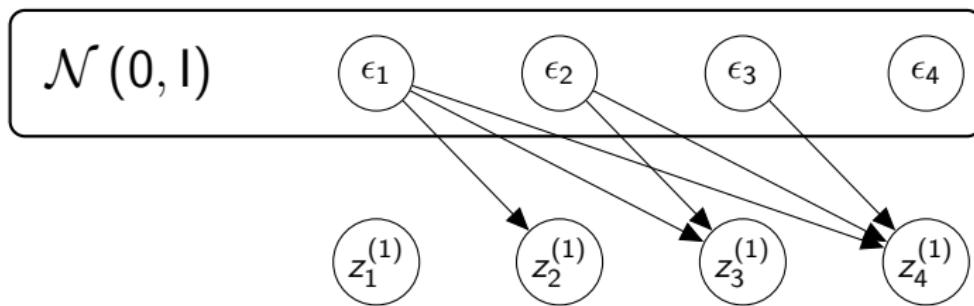
# Designing a Transformation

We are again going to use a MADE to predict parameters. However, this time we will use it in the other direction.

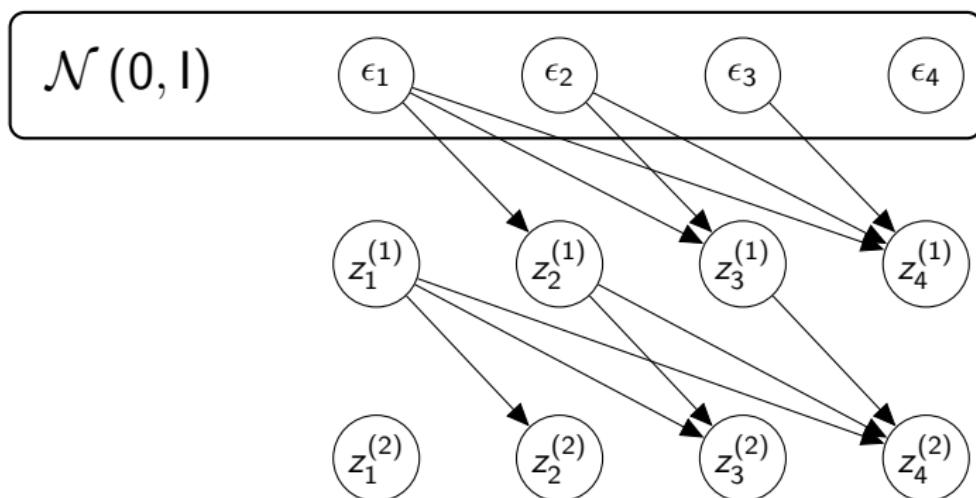
# Designing a Transformation

 $\mathcal{N}(0, I)$  $\epsilon_1$  $\epsilon_2$  $\epsilon_3$  $\epsilon_4$

# Designing a Transformation



# Designing a Transformation



# Designing a Transformation

We use a MADE  $f_\lambda$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = f_\lambda(\epsilon_{<j})$ . Then we apply the first transformation.

$$z_j^{(1)} = h_1(\epsilon)_j = \mu(\epsilon_{<j}) + \sigma(\epsilon_{<j})\epsilon_j$$

# Designing a Transformation

We use a MADE  $f_\lambda$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = f_\lambda(\epsilon_{<j})$ . Then we apply the first transformation.

$$z_j^{(1)} = h_1(\epsilon)_j = \mu(\epsilon_{<j}) + \sigma(\epsilon_{<j})\epsilon_j$$

$$z^{(1)} = h_1(\epsilon) = \mu + \sigma\epsilon$$

# Designing a Transformation

We use a MADE  $f_\lambda$  to predict the parameters of the first transformation:  $[\mu_j \ \sigma_j] = f_\lambda(\epsilon_{<j})$ . Then we apply the first transformation.

$$z_j^{(1)} = h_1(\epsilon)_j = \mu(\epsilon_{<j}) + \sigma(\epsilon_{<j})\epsilon_j$$

$$z^{(1)} = h_1(\epsilon) = \mu + \sigma\epsilon$$

$$J_{h_1}(\epsilon) = \mathbb{I}\sigma + J_\mu(\epsilon) + J_{\sigma\epsilon}(\epsilon)$$

# Designing a Transformation

## Simple Jacobian Determinant

$$|\det J_{h_1}(\epsilon)| = \prod_{j=1}^M \sigma_j$$

In practice we work with the log-likelihood.

$$\log |\det J_{h_1}(\epsilon)| = \sum_{j=1}^M \log \sigma_j$$

## 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) \left| \det J_{h_1}(z^{(1)}) \right| \left| \det J_{h_2}(z^{(2)}) \right| \\ &= q(h_1(h_2(z^{(2)}))) \left| \det J_{h_1}(h_2(z^{(2)})) \right| \left| \det J_{h_2}(z^{(2)}) \right| \end{aligned}$$

## 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) \left| \det J_{h_1}(z^{(1)}) \right| \left| \det J_{h_2}(z^{(2)}) \right| \\ &= q(h_1(h_2(z^{(2)}))) \left| \det J_{h_1}(h_2(z^{(2)})) \right| \left| \det J_{h_2}(z^{(2)}) \right| \end{aligned}$$

$$\log q(z^{(2)}) = \log q(h_1(h_2(z^{(2)}))) + \sum_{j=1}^M \log \sigma_j^{(1)} + \sum_{j=1}^M \log \sigma_j^{(2)}$$

## 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) \left| \det J_{h_1}(z^{(1)}) \right| \left| \det J_{h_2}(z^{(2)}) \right| \\ &= q(h_1(h_2(z^{(2)}))) \left| \det J_{h_1}(h_2(z^{(2)})) \right| \left| \det J_{h_2}(z^{(2)}) \right| \end{aligned}$$

$$\log q(z^{(2)}) = \log q(h_1(h_2(z^{(2)}))) + \sum_{j=1}^M \log \sigma_j^{(1)} + \sum_{j=1}^M \log \sigma_j^{(2)}$$

$$z^{(1)} = \mu^{(1)} + \sigma^{(1)}\epsilon \text{ where } [\mu^{(1)}, \sigma^{(1)}] = f_\lambda^{(1)}(\epsilon)$$

## 2-step Flow

$$\begin{aligned} q(z^{(2)}) &= q(\epsilon) \left| \det J_{h_1}(z^{(1)}) \right| \left| \det J_{h_2}(z^{(2)}) \right| \\ &= q(h_1(h_2(z^{(2)}))) \left| \det J_{h_1}(h_2(z^{(2)})) \right| \left| \det J_{h_2}(z^{(2)}) \right| \end{aligned}$$

$$\log q(z^{(2)}) = \log q(h_1(h_2(z^{(2)}))) + \sum_{j=1}^M \log \sigma_j^{(1)} + \sum_{j=1}^M \log \sigma_j^{(2)}$$

$$z^{(1)} = \mu^{(1)} + \sigma^{(1)} \epsilon \text{ where } [\mu^{(1)}, \sigma^{(1)}] = f_\lambda^{(1)}(\epsilon)$$

$$z^{(2)} = \mu^{(2)} + \sigma^{(2)} z^{(1)} \text{ where } [\mu^{(2)}, \sigma^{(2)}] = f_\lambda^{(2)}(z^{(1)})$$

# ELBO

$$\text{ELBO} = \mathbb{E}_{q(z|\lambda)} [\log p(x|z)] - \text{KL}(q(z^{(2)}|\lambda) || p(z^{(2)})) =$$

# ELBO

$$\text{ELBO} = \mathbb{E}_{q(z|\lambda)} [\log p(x|z)] - \text{KL}(q(z^{(2)}|\lambda) || p(z^{(2)})) = \\ \mathbb{E}_{q(z|\lambda)} [\log p(x|z)] - \text{KL}(q(\epsilon) |\det J_h(z^{(2)}) | || p(z))$$

# ELBO

## KL-term

$$\text{KL} \left( q(\epsilon) \left| \det J_h(z^{(2)}) \right| \parallel p(z) \right) =$$

# ELBO

## KL-term

$$\text{KL} (q(\epsilon) \left| \det J_h(z^{(2)}) \right| \parallel p(z)) = \\ \mathbb{E}_{q(z^{(2)}|\lambda))} \left[ \frac{q(\epsilon) \left| \det J_h(z^{(2)}) \right|}{p(z^{(2)})} \right]$$

# ELBO

## KL-term

$$\begin{aligned} \text{KL} (q(\epsilon) \mid\mid p(z)) &= \\ \mathbb{E}_{q(z^{(2)} | \lambda))} \left[ \frac{q(\epsilon) \mid \det J_h(z^{(2)}) \mid}{p(z^{(2)})} \right] &\stackrel{\text{MC}}{\approx} \frac{1}{S} \sum_{s=1}^S \frac{q(\epsilon) \mid \det J_h(z^{(2,s)}) \mid}{p(z^{(2,s)})} \end{aligned}$$

# ELBO

## KL-term

$$\begin{aligned} \text{KL} (q(\epsilon) \mid\mid p(z)) &= \\ \mathbb{E}_{q(z^{(2)} | \lambda))} \left[ \frac{q(\epsilon) \mid \det J_h(z^{(2)}) \mid}{p(z^{(2)})} \right] &\stackrel{\text{MC}}{\approx} \frac{1}{S} \sum_{s=1}^S \frac{q(\epsilon) \mid \det J_h(z^{(2,s)}) \mid}{p(z^{(2,s)})} \end{aligned}$$

## Jacobian

# ELBO

## KL-term

$$\begin{aligned} \text{KL} (q(\epsilon) \mid\mid p(z)) &= \\ \mathbb{E}_{q(z^{(2)} | \lambda))} \left[ \frac{q(\epsilon) \mid\det J_h(z^{(2)})\mid}{p(z^{(2)})} \right] &\stackrel{\text{MC}}{\approx} \frac{1}{S} \sum_{s=1}^S \frac{q(\epsilon) \mid\det J_h(z^{(2,s)})\mid}{p(z^{(2,s)})} \end{aligned}$$

## Jacobian

$$\mid\det J_h(z^{(2,s)})\mid = \sum_{j=1}^M \log \sigma_j^{(1)} + \sum_{j=1}^M \log \sigma_j^{(2)}$$

# Other Applications of Normalizing Flows

- ▶ As a prior

# Other Applications of Normalizing Flows

- ▶ As a prior
- ▶ Modeling of dynamic systems

# Summary

- ▶ NFs model arbitrary continuous distributions
- ▶ They allow for density computation
- ▶ Need to have simple Jacobian
- ▶ Depending on direction, they are good at either sampling or density computation (not both)

# References I

Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle. Made: Masked autoencoder for distribution estimation. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, pages 881–889, 2015.

Chin-Wei Huang, David Krueger, Alexandre Lacoste, and Aaron C. Courville. Neural autoregressive flows. *CoRR*, 2018. URL <http://arxiv.org/abs/1804.00779>.

## References II

- Diederik P. Kingma, Tim Salimans, and Max Welling. Improving variational inference with inverse autoregressive flow. 2016. URL <http://arxiv.org/abs/1606.04934>.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *ICML*, pages 1530–1538, 2015. URL <http://proceedings.mlr.press/v37/rezende15.pdf>.