

Advanced Topics in Machine Learning (deep generative modelling)

Lecture 4: Models with tractable exact density



Nikolay Malkin

3 February 2026

Outline of Lecture 4

Autoregressive models and normalising flows:

- ▶ Review
 - ▶ Warmup: Autoregressive models
- ▶ Noise outsourcing and pushforward measures
- ▶ Normalising flows

- ▶ Review

- ▶ Warmup: Autoregressive models

- ▶ Noise outsourcing and pushforward measures

- ▶ Normalising flows

Review: Generative modelling as divergence minimisation

Setting:

- ▶ We have a **data distribution** π_{data} over \mathbb{R}^d (from which we can sample, but we do not know its density function)
 - ▶ It could be the empirical distribution of a dataset
- ▶ We have a class of **model distributions** $\{\pi_{\theta}\}$ (with densities p_{θ})
 - ▶ θ are the parameters of the model (e.g., neural network weights, Gaussian mixture parameters)
 - ▶ Note that we do not necessarily know the density functions p_{θ}
- ▶ We seek θ such that π_{θ} approximates π_{data} well:

$$\theta^* = \arg \min_{\theta} D(\pi_{\theta}, \pi_{\text{data}})$$

Review: Generative modelling as divergence minimisation

The KL divergence

$$\text{KL}(p\|q) = \mathbb{E}_{X \sim p} \left[\log \frac{p(X)}{q(X)} \right] = \int_{\mathbb{R}^d} p(x) \log \frac{p(x)}{q(x)} dx$$

- ▶ Minimising $\text{KL}(\pi_{\text{data}}\|\pi_{\theta}) \equiv$ maximising sample log-likelihood $\mathbb{E}_{X \sim \pi_{\text{data}}} [\log p_{\theta}(X)]$
- ▶ Algorithm to fit θ using stochastic gradient descent:
 - ▶ Sample a minibatch $x_1, \dots, x_m \sim \pi_{\text{data}}$
 - ▶ Compute gradient estimate:

$$g = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} [-\log p_{\theta}(x_i)]$$

- ▶ Update parameters: $\theta \leftarrow \theta - \eta g$ (or using your favourite optimiser)

Autoregressive models

Simplest models with tractable exact density: **autoregressive models**

- ▶ Factor a density $p(x_1, \dots, x_d)$ as

$$p(x_1, \dots, x_d) = p(x_1)p(x_2|x_1) \cdots p(x_d|x_1, \dots, x_{d-1})$$

- ▶ Model each factor using a neural network

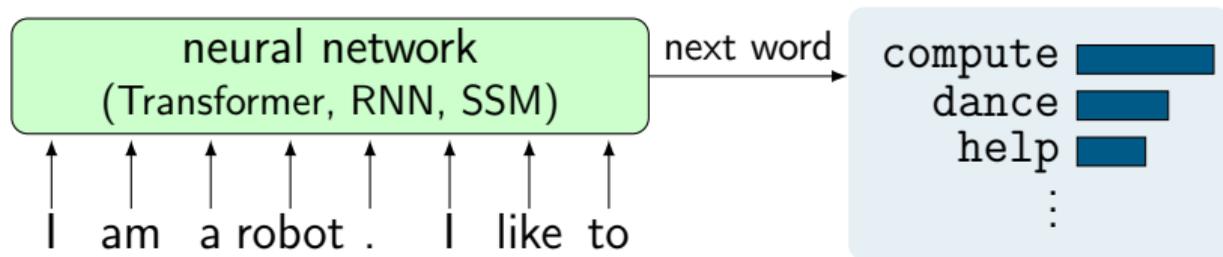
Autoregressive models

Simplest models with tractable exact density: **autoregressive models**

- ▶ Factor a density $p(x_1, \dots, x_d)$ as

$$p(x_1, \dots, x_d) = p(x_1)p(x_2|x_1) \cdots p(x_d|x_1, \dots, x_{d-1})$$

- ▶ Model each factor using a neural network



- ▶ Can also randomise order (related to discrete diffusion models)

Autoregressive models

Simplest models with tractable exact density: **autoregressive models**

- ▶ Factor a density $p(x_1, \dots, x_d)$ as

$$p(x_1, \dots, x_d) = p(x_1)p(x_2|x_1) \cdots p(x_d|x_1, \dots, x_{d-1})$$

- ▶ Model each factor using a neural network
- ▶ Common in language modelling, but also possible for images



MADE [Germain et al., 2015]



PixelRNN
[van den Oord et al., 2016]

- ▶ Can also randomise order (related to discrete diffusion models)

- ▶ Review

 - ▶ Warmup: Autoregressive models

- ▶ Noise outsourcing and pushforward measures

- ▶ Normalising flows

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$. $f_{\#}\pi = \delta_{35}$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$. $f_{\#}\pi = \delta_{35}$.
- ▶ $\pi = \mathcal{N}(0, 1)$, $f(x) = 2x + 1$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$. $f_{\#}\pi = \delta_{35}$.
- ▶ $\pi = \mathcal{N}(0, 1)$, $f(x) = 2x + 1$. $f_{\#}\pi = \mathcal{N}(1, 4)$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$. $f_{\#}\pi = \delta_{35}$.
- ▶ $\pi = \mathcal{N}(0, 1)$, $f(x) = 2x + 1$. $f_{\#}\pi = \mathcal{N}(1, 4)$.
- ▶ $\pi = \text{Unif}(0, 2\pi)$, $f(x) = (\cos x, \sin x)$.

Pushforward measure

If π is a probability distribution over \mathbb{R}^d and $f : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a **measurable** function, then the **pushforward distribution** $f_{\#}\pi$ is a distribution over $\mathbb{R}^{d'}$ that is sampled by the procedure:

- ▶ Sample $x \sim \pi$
- ▶ Output $f(x)$

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$, $f(x) = 2x + 1$. $f_{\#}\pi = \text{Unif}([1, 3])$.
- ▶ $\pi = \delta_{17}$, $f(x) = 2x + 1$. $f_{\#}\pi = \delta_{35}$.
- ▶ $\pi = \mathcal{N}(0, 1)$, $f(x) = 2x + 1$. $f_{\#}\pi = \mathcal{N}(1, 4)$.
- ▶ $\pi = \text{Unif}(0, 2\pi)$, $f(x) = (\cos x, \sin x)$. $f_{\#}\pi$ is uniform on the circle.

Short review of variable change formula

A special case:

- ▶ π has a density $p : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- ▶ $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is **invertible** and **differentiable**, with differentiable inverse f^{-1} (a **diffeomorphism**)

Short review of variable change formula

A special case:

- ▶ π has a density $p : \mathbb{R}^d \rightarrow \mathbb{R}_{\geq 0}$
- ▶ $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is **invertible** and **differentiable**, with differentiable inverse f^{-1} (a **diffeomorphism**)

Change of variables formula: if $y = f(x)$, then

density of the pushforward at y

$$\begin{aligned} \overbrace{(f_{\#}p)(y)} &= p(x) \left| \det \left(\frac{\partial f^{-1}}{\partial y}(y) \right) \right| \\ &= p(x) \left[\det \left(\frac{\partial f}{\partial x}(x) \right) \right]^{-1} \end{aligned}$$

determinant of Jacobian of f at $x = f^{-1}(y)$

Short review of variable change formula

density of the pushforward at y

$$\underbrace{(f_{\#}p)(y)} = p(x) \left| \det \left(\frac{\partial f^{-1}}{\partial y}(y) \right) \right| = p(x) \underbrace{\left| \det \left(\frac{\partial f}{\partial x}(x) \right) \right|}^{-1}$$

determinant of Jacobian of f at $x = f^{-1}(y)$

Why?

► Linear approximation:

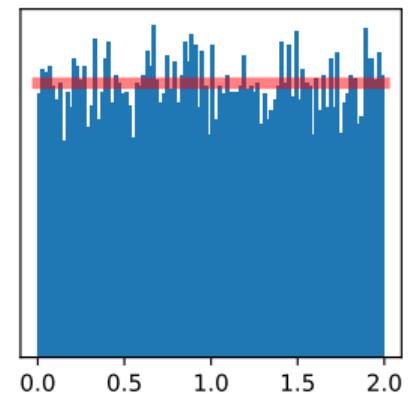
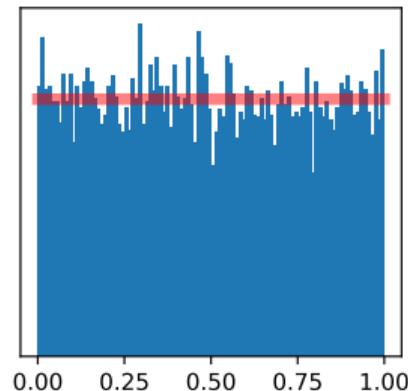
$$f(x + \Delta x) = f(x) + \underbrace{\frac{\partial f}{\partial x}(x)}_{d \times d \text{ matrix}} \underbrace{\Delta x}_{\text{vector in } \mathbb{R}^d} + O(\|\Delta x\|^2)$$

► Volumes are scaled by the determinant under a linear transformation

Short review of variable change formula

Examples:

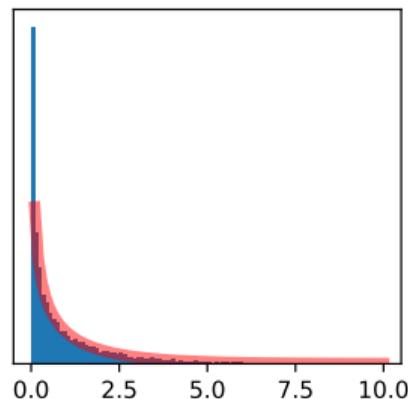
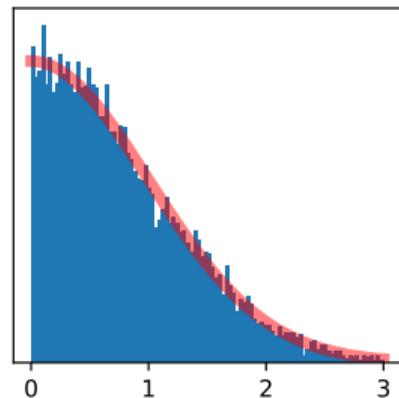
- ▶ $\pi = \text{Unif}([0, 1])$ ($p(x) = 1$ for $x \in [0, 1]$),
 $f(x) = 2x$
 - ▶ $|f'(x)| = 2$, so $(f_{\#}\pi)(y) = \frac{1}{2}$ for $y \in [0, 2]$



Short review of variable change formula

Examples:

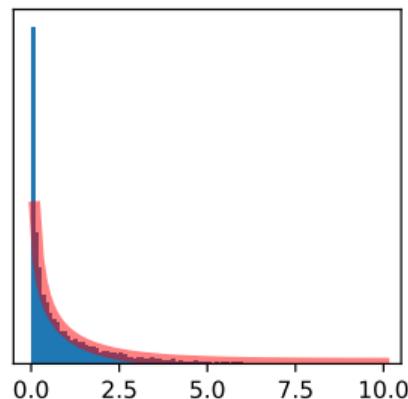
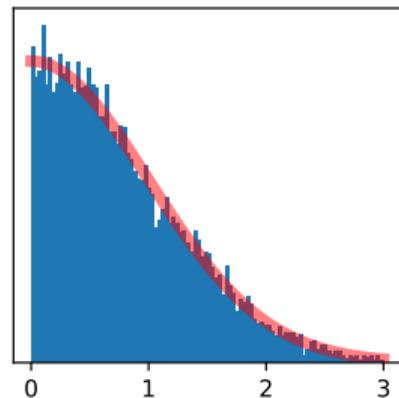
- ▶ $\pi = \text{Unif}([0, 1])$ ($p(x) = 1$ for $x \in [0, 1]$),
 $f(x) = 2x$
 - ▶ $|f'(x)| = 2$, so $(f_{\#}\pi)(y) = \frac{1}{2}$ for $y \in [0, 2]$
- ▶ $\pi = \mathcal{N}(0, 1)$ restricted to $[0, \infty)$
($p(x) = \frac{2}{\sqrt{2\pi}} \exp(\frac{-x^2}{2})$), $f(x) = x^2$
 - ▶ $|f'(x)| = 2|x|$, so $(f_{\#}\pi)(y) = \frac{1}{\sqrt{2\pi y}} \exp(\frac{-y}{2})$ for $y \geq 0$



Short review of variable change formula

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$ ($p(x) = 1$ for $x \in [0, 1]$),
 $f(x) = 2x$
 - ▶ $|f'(x)| = 2$, so $(f_{\#}\pi)(y) = \frac{1}{2}$ for $y \in [0, 2]$
- ▶ $\pi = \mathcal{N}(0, 1)$ restricted to $[0, \infty)$
($p(x) = \frac{2}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$), $f(x) = x^2$
 - ▶ $|f'(x)| = 2|x|$, so $(f_{\#}\pi)(y) = \frac{1}{\sqrt{2\pi y}} \exp(-\frac{y}{2})$ for $y \geq 0$
(χ^2 distribution with one degree of freedom).



Short review of variable change formula

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$ ($p(x) = 1$ for $x \in [0, 1]$),
 $f(x) = 2x$
 - ▶ $|f'(x)| = 2$, so $(f_{\#}\pi)(y) = \frac{1}{2}$ for $y \in [0, 2]$
- ▶ $\pi = \mathcal{N}(0, 1)$ restricted to $[0, \infty)$
($p(x) = \frac{2}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$), $f(x) = x^2$
 - ▶ $|f'(x)| = 2|x|$, so $(f_{\#}\pi)(y) = \frac{1}{\sqrt{2\pi y}} \exp(-\frac{y}{2})$ for $y \geq 0$
(χ^2 distribution with one degree of freedom).
- ▶ $\pi = \text{Unif}([0, 1] \times [0, 2\pi])$, $f(r, \theta) = (r \cos \theta, r \sin \theta)$

Short review of variable change formula

Examples:

- ▶ $\pi = \text{Unif}([0, 1])$ ($p(x) = 1$ for $x \in [0, 1]$),
 $f(x) = 2x$

- ▶ $|f'(x)| = 2$, so $(f_{\#}\pi)(y) = \frac{1}{2}$ for $y \in [0, 2]$

- ▶ $\pi = \mathcal{N}(0, 1)$ restricted to $[0, \infty)$

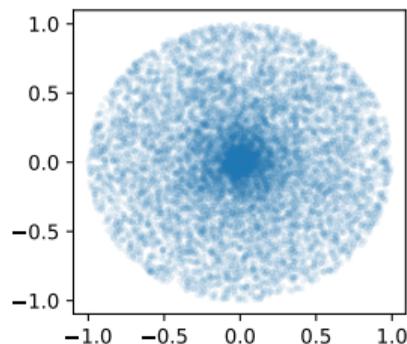
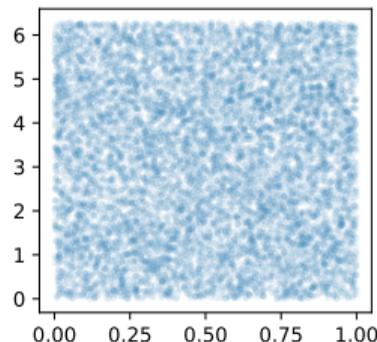
- ($p(x) = \frac{2}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$), $f(x) = x^2$

- ▶ $|f'(x)| = 2|x|$, so $(f_{\#}\pi)(y) = \frac{1}{\sqrt{2\pi y}} \exp(-\frac{y}{2})$ for $y \geq 0$
(χ^2 distribution with one degree of freedom).

- ▶ $\pi = \text{Unif}([0, 1] \times [0, 2\pi])$, $f(r, \theta) = (r \cos \theta, r \sin \theta)$

$$\det \left(\frac{\partial f}{\partial(r, \theta)} \right) = \det \begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix} = r$$

- ▶ $(f_{\#}\pi)(x, y) = \frac{1}{2\pi r}$, where $r = \sqrt{x^2 + y^2}$, for $r \in [0, 1]$



Noise outsourcing

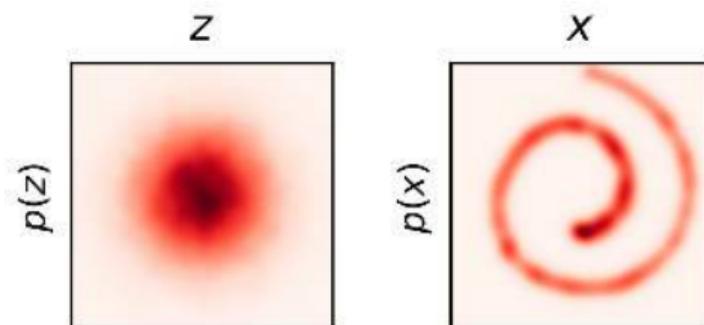
Noise outsourcing lemma: Any distribution π over \mathbb{R}^d can be represented as the pushforward of a simple distribution (e.g., $\mathcal{N}(0, I_d)$ or $\text{Unif}([0, 1]^d)$) by some function

- ▶ Not very useful, does not tell us how to find such a function or guarantee invertibility and differentiability

Noise outsourcing

Noise outsourcing lemma: Any distribution π over \mathbb{R}^d can be represented as the pushforward of a simple distribution (e.g., $\mathcal{N}(0, I_d)$ or $\text{Unif}([0, 1]^d)$) by some function

- ▶ Not very useful, does not tell us how to find such a function or guarantee invertibility and differentiability
- ▶ But if π has a density **and satisfies other basic conditions**, then such a smooth invertible function **does** exist (related to Brenier's theorem, optimal transport)



Noise outsourcing

Noise outsourcing lemma: Any distribution π over \mathbb{R}^d can be represented as the pushforward of a simple distribution (e.g., $\mathcal{N}(0, I_d)$ or $\text{Unif}([0, 1]^d)$) by some function

- ▶ Not very useful, does not tell us how to find such a function or guarantee invertibility and differentiability
- ▶ But if π has a density and satisfies other basic conditions, then such a smooth invertible function **does** exist (related to Brenier's theorem, optimal transport)

Implication for generative modelling: To approximate π_{data} :

- ▶ Parametrise a class of smooth invertible functions $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are easy to invert and differentiate
- ▶ Define $\pi_\theta = (f_\theta)_\# \pi_{\text{base}}$, where π_{base} is a simple distribution (e.g., $\mathcal{N}(0, I_d)$)
- ▶ Fit θ by maximising data log-likelihood (given by variable change formula)

- ▶ Review

 - ▶ Warmup: Autoregressive models

- ▶ Noise outsourcing and pushforward measures

- ▶ Normalising flows

Normalising flows

Parametrise a class of smooth invertible functions $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are easy to invert and differentiate

Normalising flows

Parametrise a class of smooth invertible functions $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are easy to invert and differentiate How do we build such functions? Compose simple invertible layers!

- ▶ If $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are diffeomorphisms, then so is $g \circ f$, and the inverse is $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$

Normalising flows

Parametrise a class of smooth invertible functions $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are easy to invert and differentiate How do we build such functions? Compose simple invertible layers!

- ▶ If $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are diffeomorphisms, then so is $g \circ f$, and the inverse is $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$
- ▶ Chain rule: if $y = f(x)$ and $z = g(y)$, then

$$\det \left(\frac{\partial z}{\partial x} \right) = \det \left(\frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \right) = \det \left(\frac{\partial z}{\partial y} \right) \cdot \det \left(\frac{\partial y}{\partial x} \right)$$

- ▶ Compose more such layers to get more complex functions

Normalising flows

Parametrise a class of smooth invertible functions $f_\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d$ that are easy to invert and differentiate How do we build such functions? Compose simple invertible layers!

- ▶ If $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ are diffeomorphisms, then so is $g \circ f$, and the inverse is $(g \circ f)^{-1} = f^{-1} \circ g^{-1}$
- ▶ Chain rule: if $y = f(x)$ and $z = g(y)$, then

$$\det \left(\frac{\partial z}{\partial x} \right) = \det \left(\frac{\partial z}{\partial y} \frac{\partial y}{\partial x} \right) = \det \left(\frac{\partial z}{\partial y} \right) \cdot \det \left(\frac{\partial y}{\partial x} \right)$$

- ▶ Compose more such layers to get more complex functions

Next: a library of simple layers

Linear layer

$f(x) = Ax + b$, where A is an invertible $d \times d$ matrix and $b \in \mathbb{R}^d$

▶ Inverse: $f^{-1}(y) = A^{-1}(y - b)$

▶ Jacobian determinant: $\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = |\det A|$

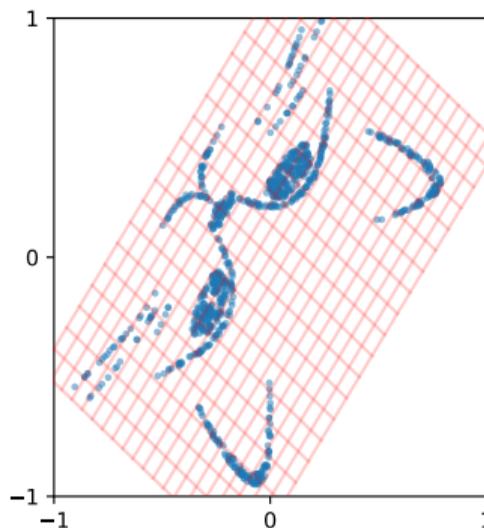
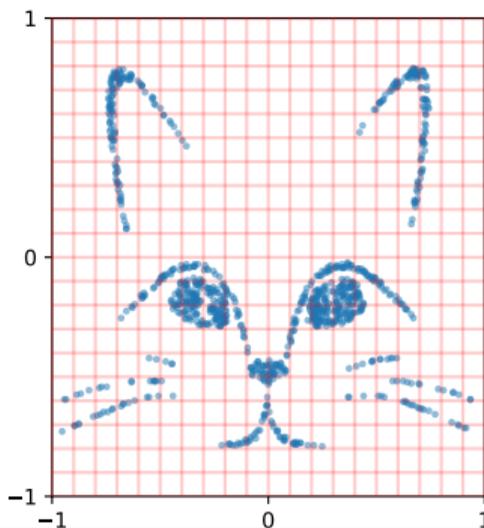
Linear layer

$f(x) = Ax + b$, where A is an invertible $d \times d$ matrix and $b \in \mathbb{R}^d$

▶ Inverse: $f^{-1}(y) = A^{-1}(y - b)$

▶ Jacobian determinant: $|\det(\frac{\partial f}{\partial x})| = |\det A|$

Can we build a model out of just linear layers?



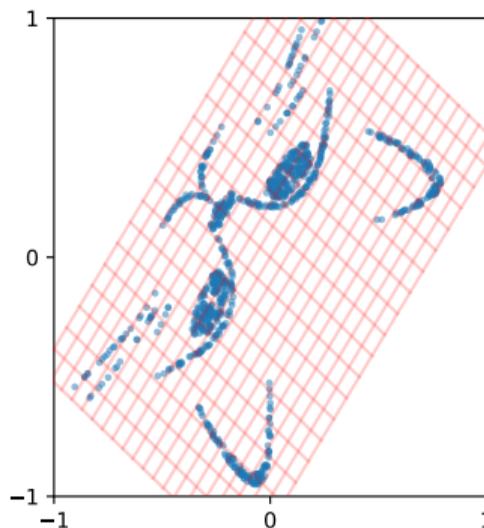
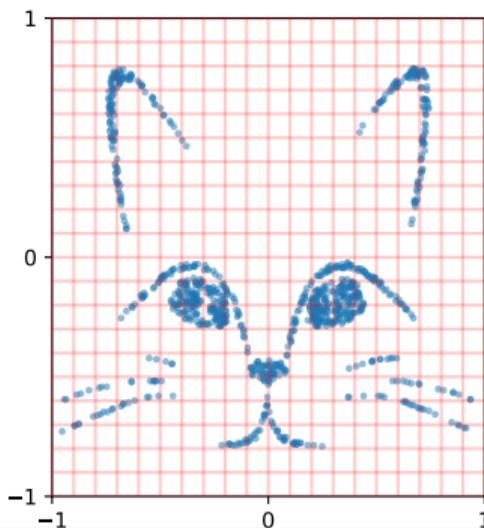
Linear layer

$f(x) = Ax + b$, where A is an invertible $d \times d$ matrix and $b \in \mathbb{R}^d$

▶ Inverse: $f^{-1}(y) = A^{-1}(y - b)$

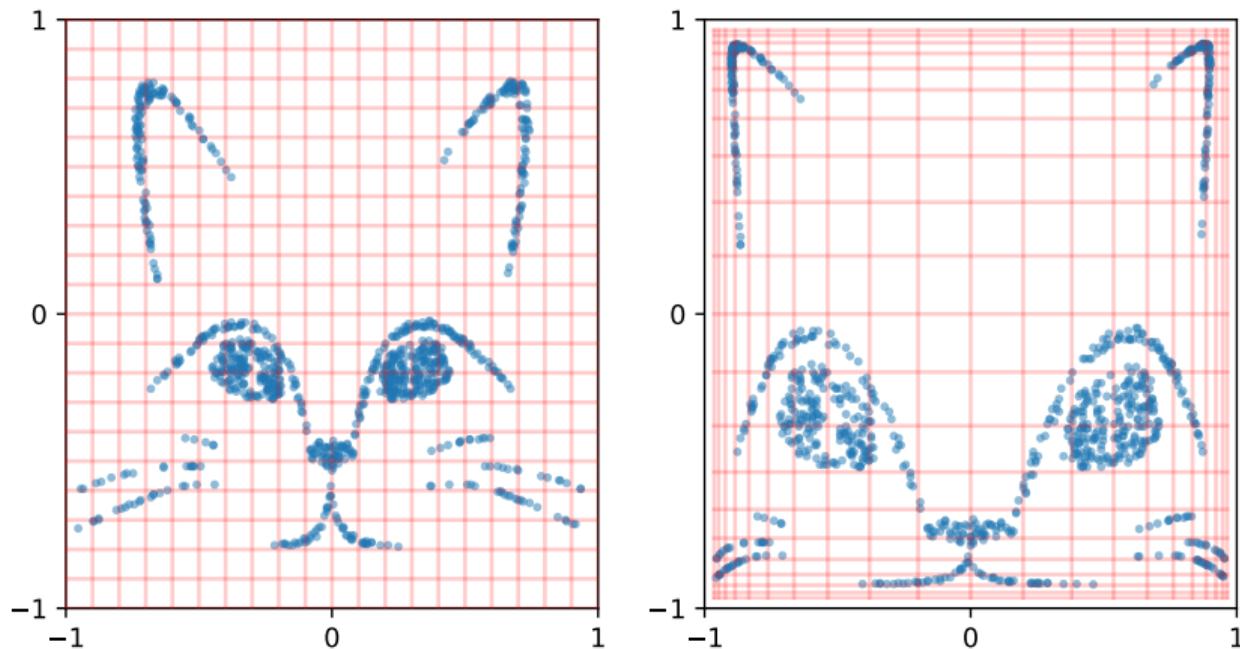
▶ Jacobian determinant: $\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = |\det A|$

Can we build a model out of just linear layers? The composition is again linear, not very expressive!



Elementwise nonlinearity layer

$f(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth invertible function (e.g., ELU, leaky ReLU, softplus, tanh)



Elementwise nonlinearity layer

$f(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth invertible function (e.g., ELU, leaky ReLU, softplus, tanh)

- ▶ Inverse: $f^{-1}(y_1, \dots, y_d) = (\sigma^{-1}(y_1), \dots, \sigma^{-1}(y_d))$
- ▶ Jacobian determinant:

$$\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = \left| \det \begin{pmatrix} \sigma'(x_1) & 0 & \cdots & 0 \\ 0 & \sigma'(x_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(x_d) \end{pmatrix} \right| = \prod_{i=1}^d |\sigma'(x_i)|$$

Elementwise nonlinearity layer

$f(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth invertible function (e.g., ELU, leaky ReLU, softplus, tanh)

▶ Inverse: $f^{-1}(y_1, \dots, y_d) = (\sigma^{-1}(y_1), \dots, \sigma^{-1}(y_d))$

▶ Jacobian determinant:

$$\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = \left| \det \begin{pmatrix} \sigma'(x_1) & 0 & \cdots & 0 \\ 0 & \sigma'(x_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(x_d) \end{pmatrix} \right| = \prod_{i=1}^d |\sigma'(x_i)|$$

Can we build a model out of just elementwise nonlinearities?

Elementwise nonlinearity layer

$f(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth invertible function (e.g., ELU, leaky ReLU, softplus, tanh)

- ▶ Inverse: $f^{-1}(y_1, \dots, y_d) = (\sigma^{-1}(y_1), \dots, \sigma^{-1}(y_d))$
- ▶ Jacobian determinant:

$$\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = \left| \det \begin{pmatrix} \sigma'(x_1) & 0 & \cdots & 0 \\ 0 & \sigma'(x_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(x_d) \end{pmatrix} \right| = \prod_{i=1}^d |\sigma'(x_i)|$$

Can we build a model out of just elementwise nonlinearities? **The composition is again elementwise, not very expressive!**
What about linear layers + elementwise nonlinearities?

Elementwise nonlinearity layer

$f(x_1, \dots, x_d) = (\sigma(x_1), \dots, \sigma(x_d))$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is a smooth invertible function (e.g., ELU, leaky ReLU, softplus, tanh)

► Inverse: $f^{-1}(y_1, \dots, y_d) = (\sigma^{-1}(y_1), \dots, \sigma^{-1}(y_d))$

► Jacobian determinant:

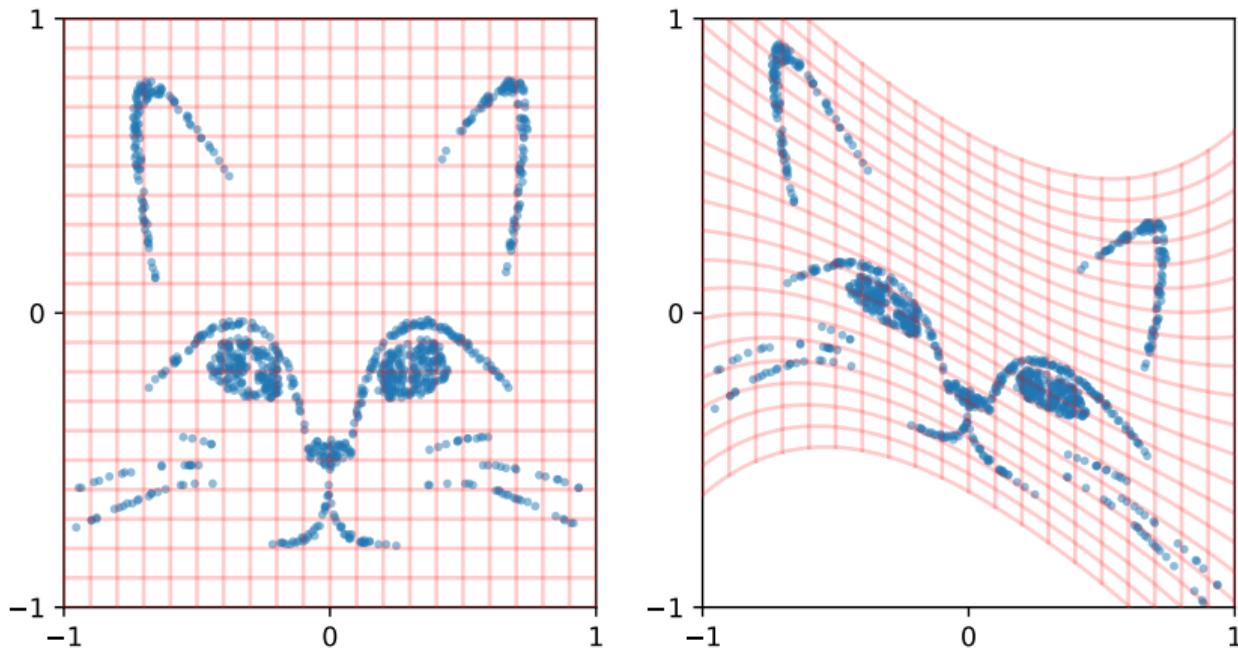
$$\left| \det \left(\frac{\partial f}{\partial x} \right) \right| = \left| \det \begin{pmatrix} \sigma'(x_1) & 0 & \cdots & 0 \\ 0 & \sigma'(x_2) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma'(x_d) \end{pmatrix} \right| = \prod_{i=1}^d |\sigma'(x_i)|$$

Can we build a model out of just elementwise nonlinearities? **The composition is again elementwise, not very expressive!**

What about linear layers + elementwise nonlinearities? **Expressive in theory, but inefficient, since dimension of all intermediate representations is d .**

Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$



Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Inverse: $f^{-1}(u, v) = (u, g_\theta^{-1}(u, v))$, where $g_\theta^{-1}(u, v) = \exp(-s_\theta(u))(v - t_\theta(u))$
- ▶ Jacobian determinant: if $(u, v) = f(x, y)$, then:

$$\left| \det \left(\frac{\partial f}{\partial(x, y)} \right) \right|$$

Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Inverse: $f^{-1}(u, v) = (u, g_\theta^{-1}(u, v))$, where $g_\theta^{-1}(u, v) = \exp(-s_\theta(u))(v - t_\theta(u))$
- ▶ Jacobian determinant: if $(u, v) = f(x, y)$, then:

$$\left| \det \left(\frac{\partial f}{\partial(x, y)} \right) \right| = \left| \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right|$$

Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Inverse: $f^{-1}(u, v) = (u, g_\theta^{-1}(u, v))$, where $g_\theta^{-1}(u, v) = \exp(-s_\theta(u))(v - t_\theta(u))$
- ▶ Jacobian determinant: if $(u, v) = f(x, y)$, then:

$$\begin{aligned} \left| \det \left(\frac{\partial f}{\partial(x, y)} \right) \right| &= \left| \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| \\ &= \left| \det \begin{pmatrix} 1 & 0 \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| \end{aligned}$$

Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Inverse: $f^{-1}(u, v) = (u, g_\theta^{-1}(u, v))$, where $g_\theta^{-1}(u, v) = \exp(-s_\theta(u))(v - t_\theta(u))$
- ▶ Jacobian determinant: if $(u, v) = f(x, y)$, then:

$$\begin{aligned} \left| \det \left(\frac{\partial f}{\partial(x, y)} \right) \right| &= \left| \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| \\ &= \left| \det \begin{pmatrix} 1 & 0 \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| = \left| \frac{\partial v}{\partial y} \right| \end{aligned}$$

Coupling layer

In two dimensions: $f(x, y) = (x, g_\theta(x, y))$, where $g : \mathbb{R}^2 \rightarrow \mathbb{R}$ is invertible in y , e.g., $g_\theta(x, y) = \exp(s_\theta(x))y + t_\theta(x)$, where $s, t : \mathbb{R} \rightarrow \mathbb{R}$

- ▶ Inverse: $f^{-1}(u, v) = (u, g_\theta^{-1}(u, v))$, where $g_\theta^{-1}(u, v) = \exp(-s_\theta(u))(v - t_\theta(u))$
- ▶ Jacobian determinant: if $(u, v) = f(x, y)$, then:

$$\begin{aligned} \left| \det \left(\frac{\partial f}{\partial(x, y)} \right) \right| &= \left| \det \begin{pmatrix} \frac{\partial u}{\partial x} & \frac{\partial u}{\partial y} \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| \\ &= \left| \det \begin{pmatrix} 1 & 0 \\ \frac{\partial v}{\partial x} & \frac{\partial v}{\partial y} \end{pmatrix} \right| = \left| \frac{\partial v}{\partial y} \right| = |\exp(s_\theta(x))| \end{aligned}$$

Coupling layer

Higher-dimensional version (transforming $x \in \mathbb{R}^d$):

- ▶ Split x in half: $x = \overbrace{x_L \oplus x_R}^{\text{concatenation}}$ with $x_L \in \mathbb{R}^{d_L}$, $x_R \in \mathbb{R}^{d_R}$ ($d_L + d_R = d$)

Coupling layer

Higher-dimensional version (transforming $x \in \mathbb{R}^d$):

- ▶ Split x in half: $x = \overbrace{x_L \oplus x_R}^{\text{concatenation}}$ with $x_L \in \mathbb{R}^{d_L}$, $x_R \in \mathbb{R}^{d_R}$ ($d_L + d_R = d$)
- ▶ Define

$$f(x_L, x_R) = \underbrace{x_L}_u \oplus \underbrace{g_\theta(x_L, x_R)}_v,$$

where $g_\theta(x_L, x_R) = \underbrace{\exp(s_\theta(x_L)) \odot x_R + t_\theta(x_L)}_{\text{elementwise product}}$ with $s_\theta, t_\theta : \mathbb{R}^{d_L} \rightarrow \mathbb{R}^{d_R}$

Coupling layer

Higher-dimensional version (transforming $x \in \mathbb{R}^d$):

- ▶ Split x in half: $x = \overbrace{x_L \oplus x_R}^{\text{concatenation}}$ with $x_L \in \mathbb{R}^{d_L}$, $x_R \in \mathbb{R}^{d_R}$ ($d_L + d_R = d$)
- ▶ Define

$$f(x_L, x_R) = \underbrace{x_L}_u \oplus \underbrace{g_\theta(x_L, x_R)}_v,$$

where $g_\theta(x_L, x_R) = \underbrace{\exp(s_\theta(x_L)) \odot x_R}_{\text{elementwise product}} + t_\theta(x_L)$ with $s_\theta, t_\theta : \mathbb{R}^{d_L} \rightarrow \mathbb{R}^{d_R}$

- ▶ Jacobian has a block structure, determinant simplifies in a similar way:

$$\left| \det \left(\frac{\partial f}{\partial (x_L, x_R)} \right) \right| = \left| \det \begin{pmatrix} I_{d_L} & 0 \\ \frac{\partial v}{\partial x_L} & \frac{\partial v}{\partial x_R} \end{pmatrix} \right| = \left| \det \left(\frac{\partial v}{\partial x_R} \right) \right| = \prod_{i=1}^{d_R} |\exp(s_\theta(x_L)_i)|$$

Coupling layer

Higher-dimensional version (transforming $x \in \mathbb{R}^d$):

- ▶ Split x in half: $x = \overbrace{x_L \oplus x_R}^{\text{concatenation}}$ with $x_L \in \mathbb{R}^{d_L}$, $x_R \in \mathbb{R}^{d_R}$ ($d_L + d_R = d$)
- ▶ Define

$$f(x_L, x_R) = \underbrace{x_L}_u \oplus \underbrace{g_\theta(x_L, x_R)}_v,$$

where $g_\theta(x_L, x_R) = \underbrace{\exp(s_\theta(x_L)) \odot x_R}_v + t_\theta(x_L)$ with $s_\theta, t_\theta : \mathbb{R}^{d_L} \rightarrow \mathbb{R}^{d_R}$

- ▶ Jacobian has a block structure, determinant simplifies in a similar way:

$$\left| \det \left(\frac{\partial f}{\partial (x_L, x_R)} \right) \right| = \left| \det \begin{pmatrix} I_{d_L} & 0 \\ \frac{\partial v}{\partial x_L} & \frac{\partial v}{\partial x_R} \end{pmatrix} \right| = \left| \det \left(\frac{\partial v}{\partial x_R} \right) \right| = \prod_{i=1}^{d_R} |\exp(s_\theta(x_L)_i)|$$

- ▶ By stacking such layers with different splits, get expressive models

Scaling to image data

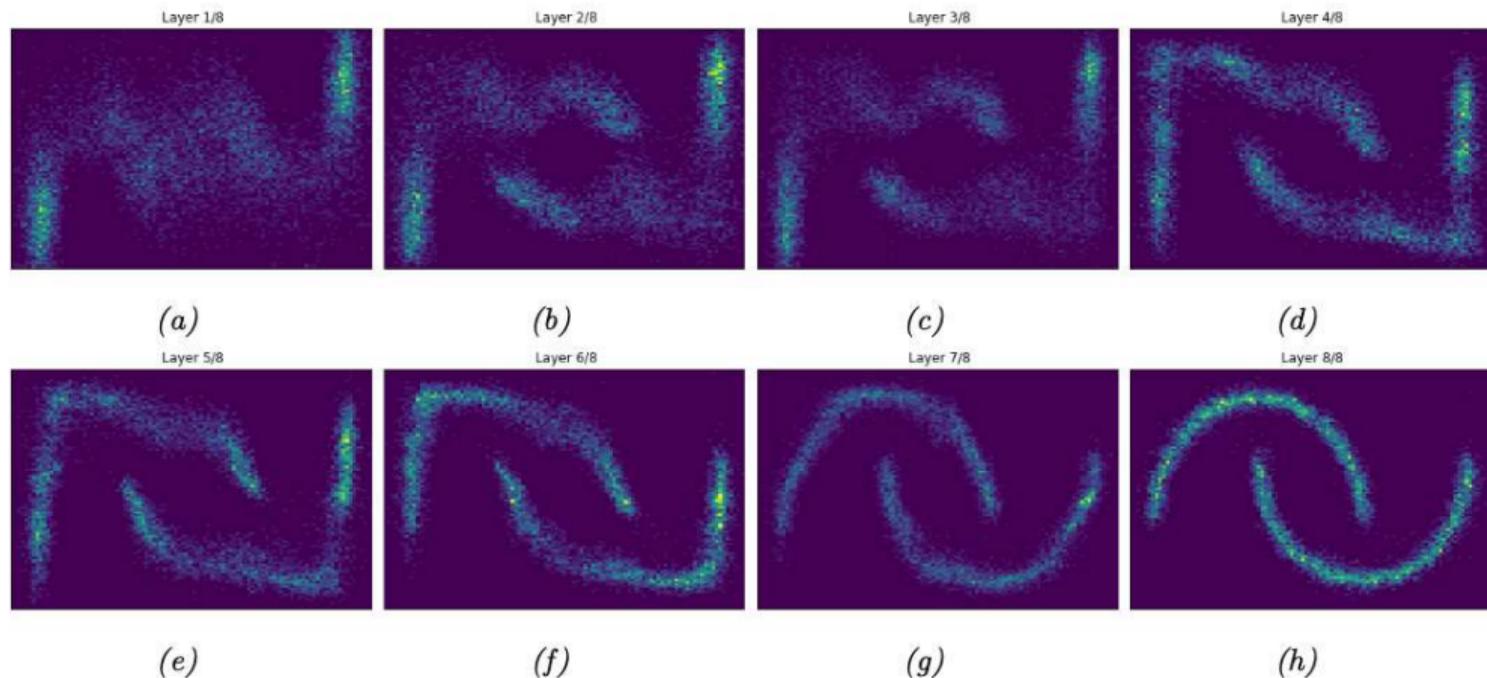
With a few domain-specific layers, normalising flows can scale to high-dimensional data (such as images)



Glow [Kingma and Dhariwal, 2018] (coupling layers with convolutions)

Some challenges of normalising flows

- ▶ Limited capacity \rightsquigarrow need many layers, mode connectivity issue



[from Murphy, §23.2]

Some challenges of normalising flows

- ▶ Limited capacity \rightsquigarrow need many layers, mode connectivity issue
- ▶ Less scalable to high-dimensional data than other models
- ▶ Requires a noise of the same dimension as the data (no latent compression)
 - ▶ If f_θ is a normalising flow and $x \sim \pi_{\text{data}}$, can we think of $z = f_\theta^{-1}(x)$ as a latent code for x ?

Some challenges of normalising flows

- ▶ Limited capacity \rightsquigarrow need many layers, mode connectivity issue
- ▶ Less scalable to high-dimensional data than other models
- ▶ Requires a noise of the same dimension as the data (no latent compression)
 - ▶ If f_θ is a normalising flow and $x \sim \pi_{\text{data}}$, can we think of $z = f_\theta^{-1}(x)$ as a latent code for x ?
 - ▶ Not really, since z has the same dimension as x and is rarely “simpler” than x
 - ▶ But we have more control with continuous normalising flows (last lectures) – invertible maps given by simulating a dynamical system

Conclusion and looking ahead

- ▶ Models with tractable exact density allow direct maximum likelihood training and evaluation
- ▶ Less scalable than other families, so fell out of favour for high-dimensional generative modelling
 - ▶ But also used as variational posteriors (in latent-variable models / for Bayesian inference)
- ▶ Next time: GANs (another pushforward model family, but trained very differently)