

# Unconstrained Monotonic Neural Networks

Accepted @ Neural Information Processing Systems 2019

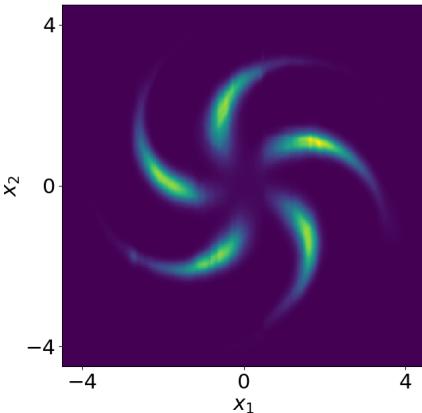
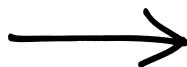
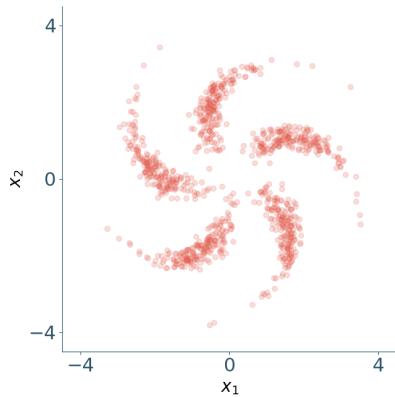
Benelearn 2019



Antoine Wehenkel

Gilles Louppe

# Density Estimation



Density estimation aims at estimating the pdf of underlying data from an iid dataset.

## Applications:

- Model uncertainty.
- Out Of Distribution detection.
- Sampling.

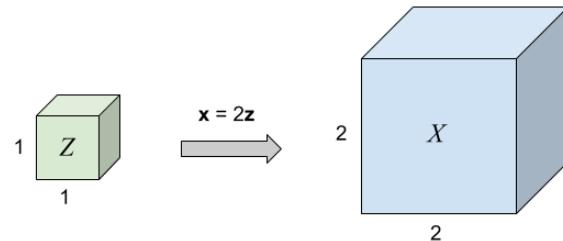
# Change of Variables Theorem (1)

Given a random variable  $\mathcal{Z}$  and a bijective function  $f$ , how does the density of  $\mathcal{X} = f(\mathcal{Z})$  behave in terms of  $p(\mathbf{z})$  and  $f$ ?

# Change of Variables Theorem (1)

Given a random variable  $\mathcal{Z}$  and a bijective function  $f$ , how does the density of  $\mathcal{X} = f(\mathcal{Z})$  behave in terms of  $p(\mathbf{z})$  and  $f$ ?

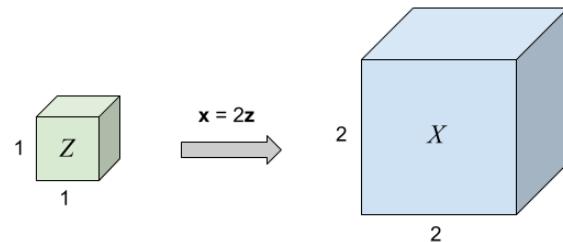
Assume  $p(z)$  is a uniformly distributed unit cube in  $\mathbb{R}^3$ , and  $\mathbf{x} = f(\mathbf{z}) = 2\mathbf{z}$ .



# Change of Variables Theorem (1)

Given a random variable  $\mathcal{Z}$  and a bijective function  $f$ , how does the density of  $\mathcal{X} = f(\mathcal{Z})$  behave in terms of  $p(\mathbf{z})$  and  $f$ ?

Assume  $p(z)$  is a uniformly distributed unit cube in  $\mathbb{R}^3$ , and  $\mathbf{x} = f(\mathbf{z}) = 2\mathbf{z}$ .



The total probability mass must be conserved, therefore

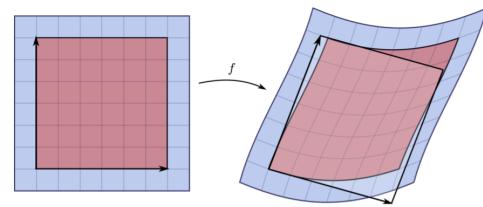
$$p(\mathbf{x}) = p(\mathbf{z}) \frac{V_z}{V_x} = p(\mathbf{z}) \frac{1}{8}, \text{ where } \frac{1}{8} = \left| \det \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 2 \end{pmatrix} \right|^{-1}$$

is the determinant of the linear transformation  $f$ .

# Change of Variables Theorem (2)

What if the transformation is non linear?

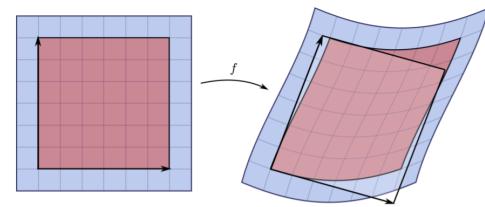
- The Jacobian  $J_f(\mathbf{z})$  of  $f(\mathbf{z}) = \mathbf{x}$  represents the infinitesimal linear transformation in the neighbourhood of  $\mathbf{z}$ .



# Change of Variables Theorem (2)

What if the transformation is non linear?

- The Jacobian  $J_f(\mathbf{z})$  of  $f(\mathbf{z}) = \mathbf{x}$  represents the infinitesimal linear transformation in the neighbourhood of  $\mathbf{z}$ .



- If the function is bijective map then the mass must be conserved locally.

Therefore, we can compute the local change of density as

$$p(\mathbf{x}) = p(\mathbf{z}) |\det J_f(\mathbf{z})|^{-1}.$$

# Change of Variables Theorem (3)

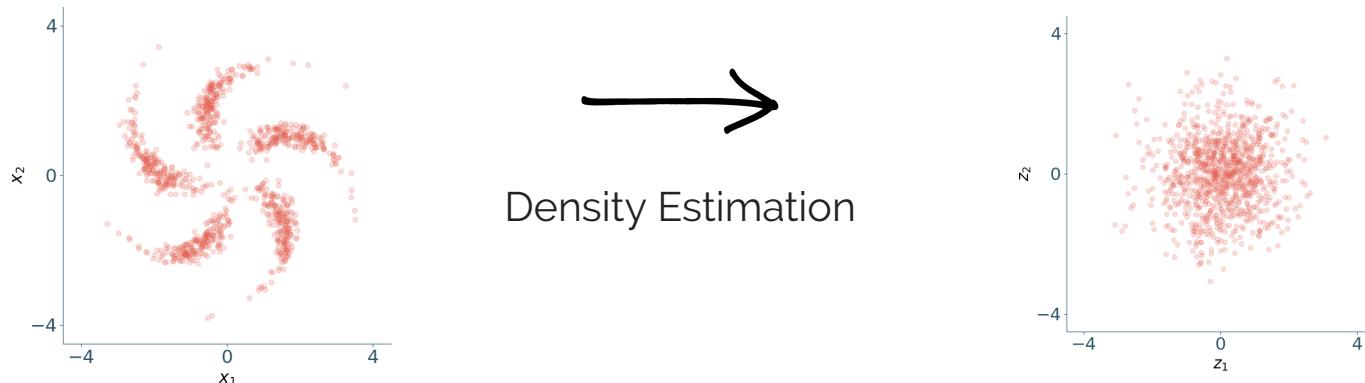
The combination of the right bijective map and any base distribution allows to represent any continuous random variable.

# Change of Variables Theorem (3)

The combination of the right bijective map and any base distribution allows to represent any continuous random variable.

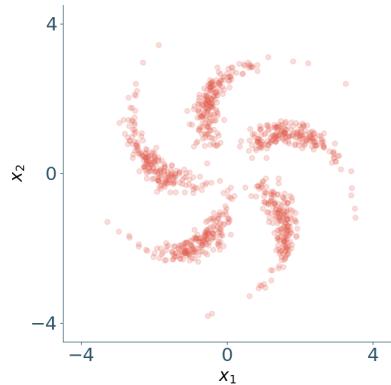
$$p(\mathbf{x}; \theta) = p(\mathbf{z} = \mathbf{g}(\mathbf{x}; \theta)) |\det J_g(\mathbf{x}; \theta)|, \quad \mathbf{g}(.; \theta) \text{ a neural network.}$$

- The bijective function takes in samples and maps them to noise.
- This process is referred as normalization if the noise distribution is normal.

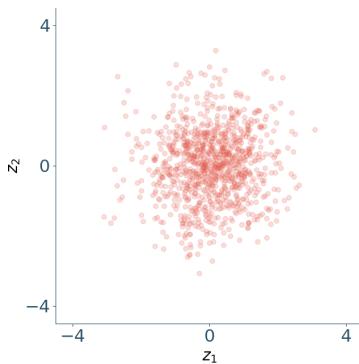


# Change of Variables Theorem (3)

The combination of the right bijective map and any base distribution allows to represent any continuous random variable.



Sampling



Once learned, the function can be inverted in order to generate samples.

# Bijectivity with Neural Nets? 🤔

- $[z_1 \dots z_d] = g([x_1 \dots x_d])$ ,  $g$  can be a NN.
- $g$  is autoregressive if it can be decomposed as:  $z_i = g_i([x_1 \dots x_i])$
- If the  $g_i$  are invertible with respect to  $x_i \forall i$ ,  $g$  is bijective.

# Bijectivity with Neural Nets? 🤔

- $[z_1 \dots z_d] = g([x_1 \dots x_d])$ ,  $g$  can be a NN.
- $g$  is autoregressive if it can be decomposed as:  $z_i = g_i([x_1 \dots x_i])$
- If the  $g_i$  are invertible with respect to  $x_i \forall i$ ,  $g$  is bijective.

The determinant of the Jacobian can be efficiently computed.

The Jacobian of an autoregressive transformation has the following form:

$$J_F(\mathbf{x}) = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} & \frac{\partial g_1}{\partial x_3} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & \frac{\partial g_2}{\partial x_3} \\ \frac{\partial g_3}{\partial x_1} & \frac{\partial g_3}{\partial x_2} & \frac{\partial g_3}{\partial x_3} \end{bmatrix} = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & 0 & 0 \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & 0 \\ \frac{\partial g_3}{\partial x_1} & \frac{\partial g_3}{\partial x_2} & \frac{\partial g_3}{\partial x_3} \end{bmatrix}.$$

# Bijectivity with Neural Nets? 🤔

- $[z_1 \dots z_d] = g([x_1 \dots x_d])$ ,  $g$  can be a NN.
- $g$  is autoregressive if it can be decomposed as:  $z_i = g_i([x_1 \dots x_i])$
- If the  $g_i$  are invertible with respect to  $x_i \forall i$ ,  $g$  is bijective.

The determinant of the Jacobian can be efficiently computed.

The Jacobian of an autoregressive transformation has the following form:

$$J_F(\mathbf{x}) = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & \frac{\partial g_1}{\partial x_2} & \frac{\partial g_1}{\partial x_3} \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & \frac{\partial g_2}{\partial x_3} \\ \frac{\partial g_3}{\partial x_1} & \frac{\partial g_3}{\partial x_2} & \frac{\partial g_3}{\partial x_3} \end{bmatrix} = \begin{bmatrix} \frac{\partial g_1}{\partial x_1} & 0 & 0 \\ \frac{\partial g_2}{\partial x_1} & \frac{\partial g_2}{\partial x_2} & 0 \\ \frac{\partial g_3}{\partial x_1} & \frac{\partial g_3}{\partial x_2} & \frac{\partial g_3}{\partial x_3} \end{bmatrix}.$$

## Chain Rule

An autoregressive density estimator learns the chain rule's factors:

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^d p(x_i | x_1, \dots, x_{i-1}).$$

# Example: Masked Autoregressive Networks

Idea: Autoregressive Networks combined with linear transformations.

- $z_1 = \sigma_1 \times x_1 + \mu_1$
- $z_i = \sigma_i(x_1, \dots, x_{i-1}) \times x_i + \mu_i(x_1, \dots, x_{i-1})$

# Example: Masked Autoregressive Networks

Idea: Autoregressive Networks combined with linear transformations.

- $z_1 = \sigma_1 \times x_1 + \mu_1$
- $z_i = \sigma_i(x_1, \dots, x_{i-1}) \times x_i + \mu_i(x_1, \dots, x_{i-1})$

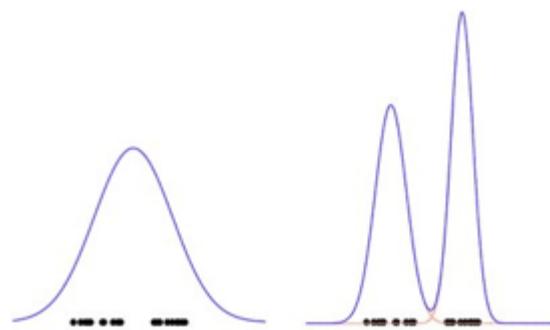
Invertible?

- $x_1 = g_1^{-1}([z_1 \dots z_d]) = g_1^{-1}(z_1) = \frac{(z_1 - \mu_1)}{\sigma_1}$
- $x_i = \frac{z_i - \mu_i([x_1 \dots x_{i-1}])}{\sigma_i([x_1 \dots x_{i-1}])}$

# Example: Masked Autoregressive Networks

Idea: Autoregressive Networks combined with linear transformations.

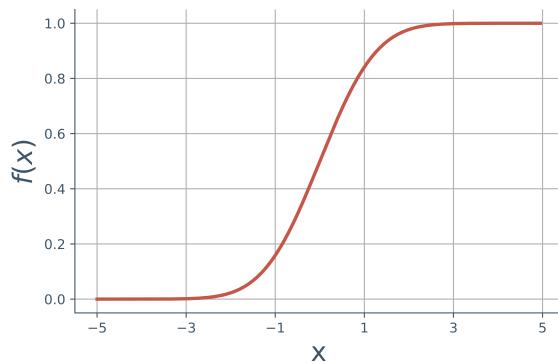
But linear transformations are not very expressive:



# Monotonic Transformations

How can we build a 1D invertible function? **Monotonic transformations.**

How can we enforce the monotonicity of a function modeled by a neural network ?



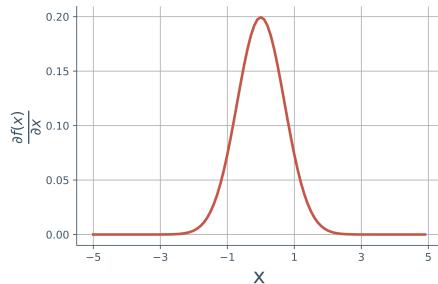
Possible solution: Strictly positive weights and monotonic activation functions.

But it restrains the architectural choice.

# Monotonic Transformations

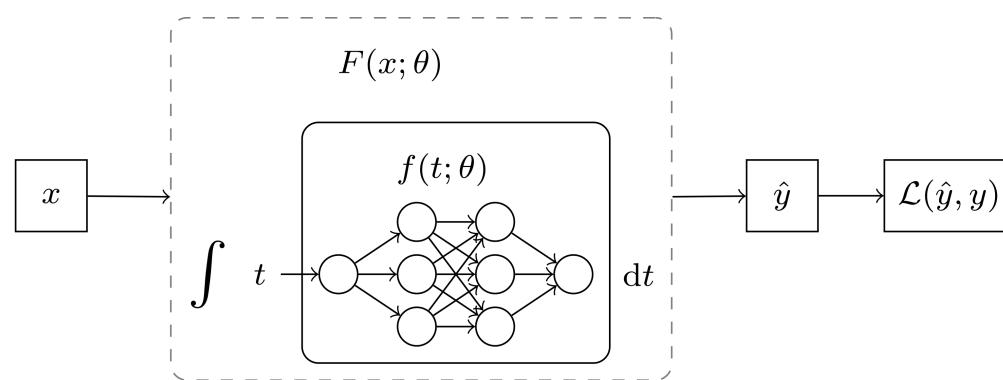
How can we enforce the monotonicity of a function modeled by a neural network ?

Our solution: To model and integrate the derivative.



The only constraint is on the output value which must be of constant sign (e.g. positive).

This can be achieved by applying an exponential on the output neuron.



# Learning of UMNN



How can we backward through the numerical integrator?

By keeping track of all the computations and building the corresponding computation graph.

# Learning of UMNN



How can we backward through the numerical integrator?

By keeping track of all the computations and building the corresponding computation graph.

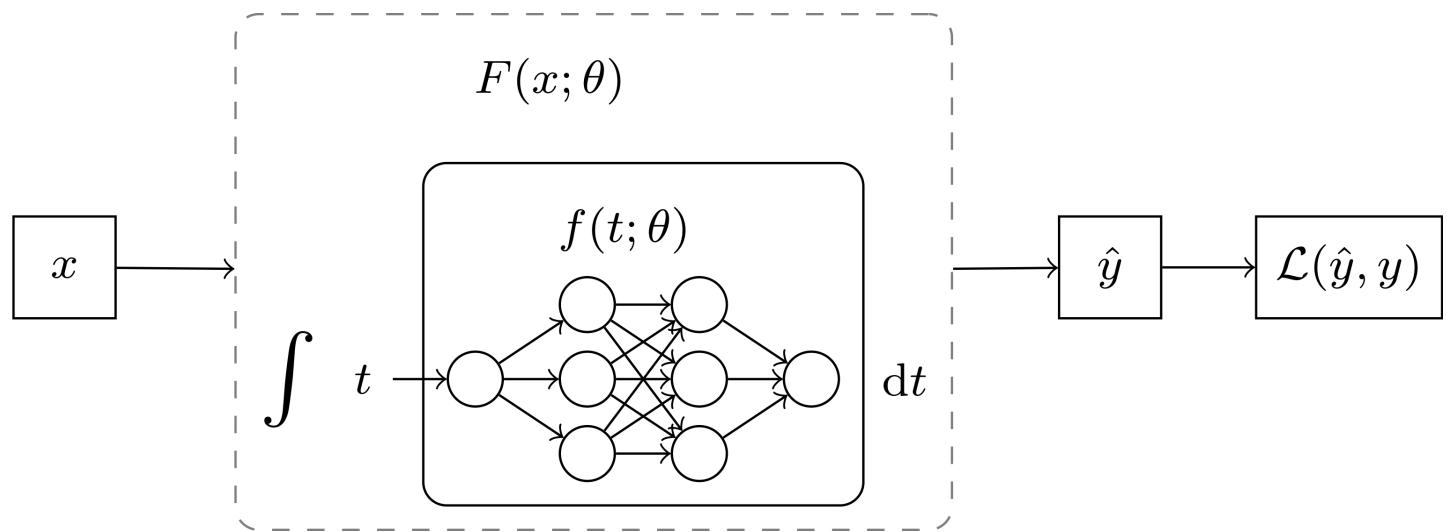


# Learning of UMNN 😠

How can we backward through the numerical integrator?

By thanking Leibnitz for his rule and applying it.

$$\frac{d}{d\omega} \left( \int_{a(\omega)}^{b(\omega)} f(t; \omega) dt \right) = f(b(\omega); \omega) \frac{d}{d\omega} b(\omega) - f(a(\omega); \omega) \frac{d}{d\omega} a(\omega) + \int_{a(\omega)}^{b(\omega)} \frac{\partial}{\partial \omega} f(t; \omega) dt$$



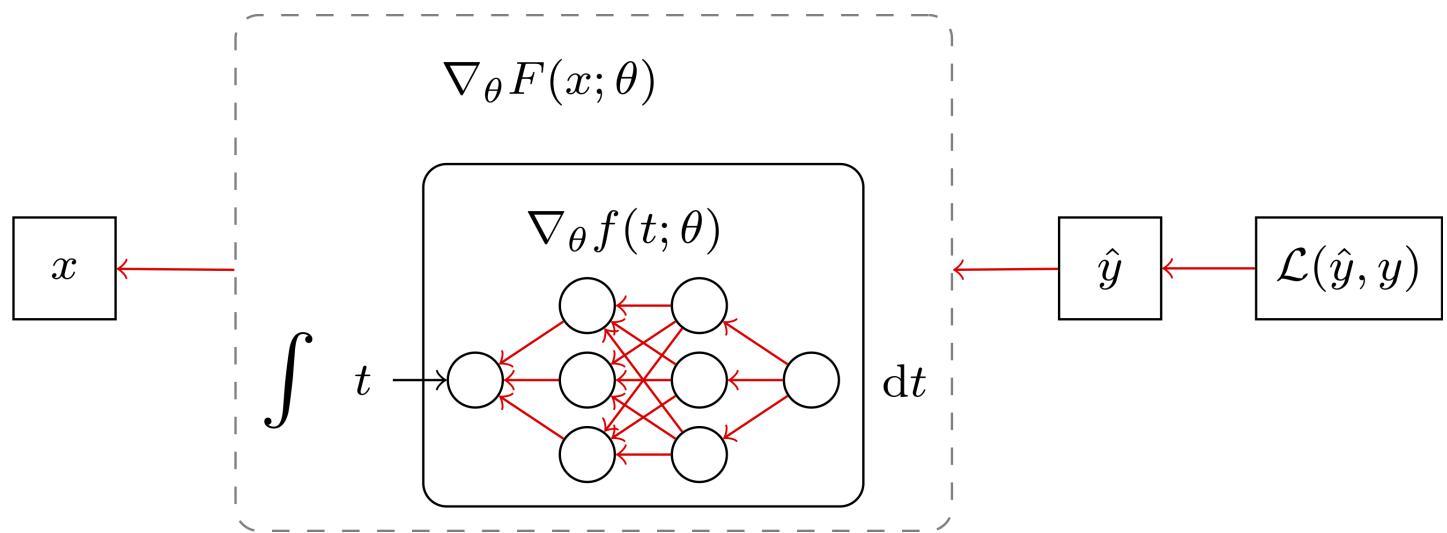
# Learning of UMNN



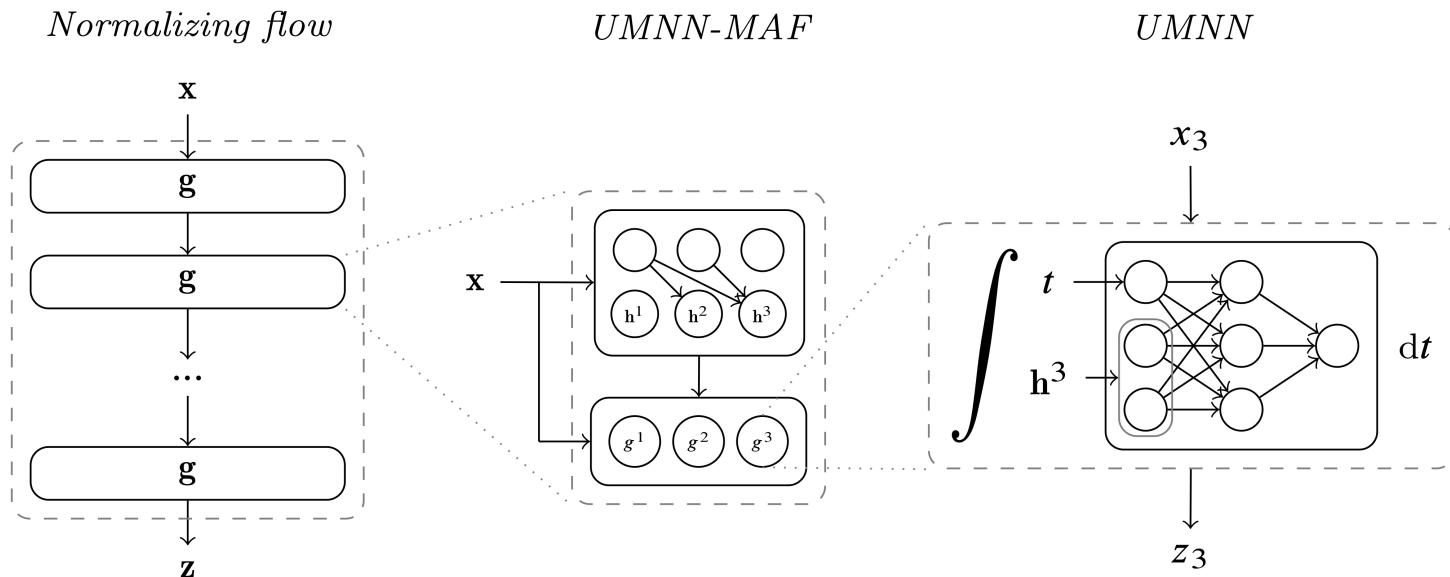
How can we backward through the numerical integrator?

By thanking Leibnitz for his rule and applying it.

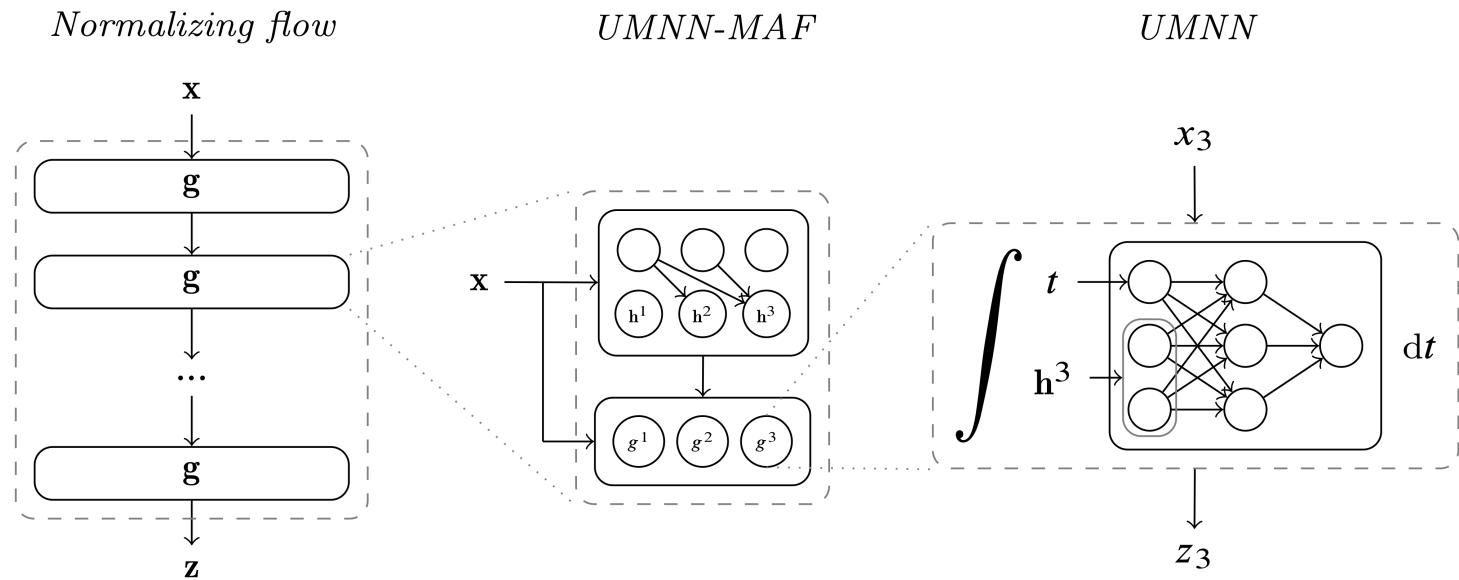
$$\frac{d}{d\omega} \left( \int_{a(\omega)}^{b(\omega)} f(t; \omega) dt \right) = f(b(\omega); \omega) \frac{d}{d\omega} b(\omega) - f(a(\omega); \omega) \frac{d}{d\omega} a(\omega) + \int_{a(\omega)}^{b(\omega)} \frac{\partial}{\partial \omega} f(t; \omega) dt$$



# Density Estimation with UMNN: UMNN-MAF



# Density Estimation with UMNN: UMNN-MAF

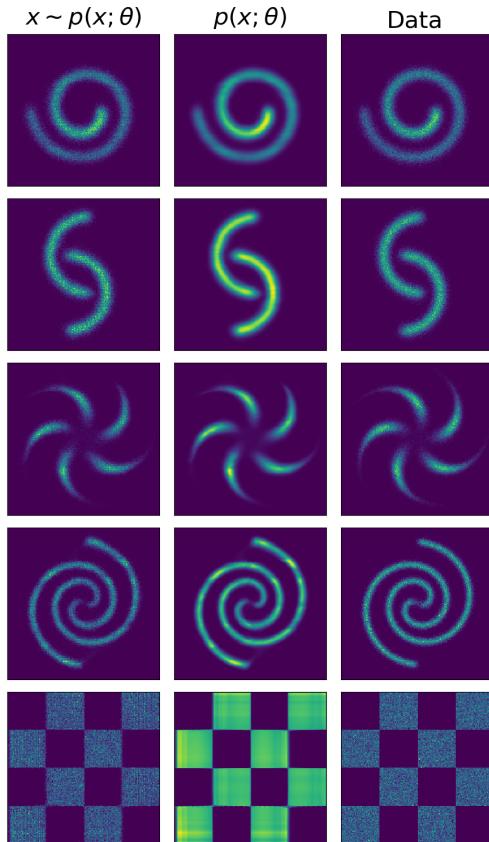


The Jacobian has the following form:

$$J_F(\mathbf{x}) = \begin{bmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_1}{\partial x_2} & \frac{\partial F_1}{\partial x_3} \\ \frac{\partial F_2}{\partial x_1} & \frac{\partial F_2}{\partial x_2} & \frac{\partial F_2}{\partial x_3} \\ \frac{\partial F_3}{\partial x_1} & \frac{\partial F_3}{\partial x_2} & \frac{\partial F_3}{\partial x_3} \end{bmatrix} = \begin{bmatrix} f_1(x_1) & 0 & 0 \\ \frac{\partial F_2}{\partial x_1} & f_2(x_2; x_1) & 0 \\ \frac{\partial F_3}{\partial x_1} & \frac{\partial F_3}{\partial x_2} & f_3(x_3; x_1, x_2) \end{bmatrix}$$

# Experimental Results

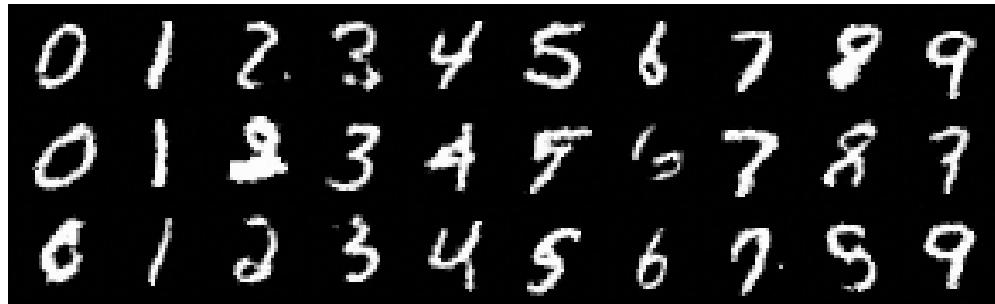
## Toy Problems



- Learnt by maximum likelihood on the train samples.
- Data dimensionality is 2.
- Multimodality.
- Discontinuity.

# Experimental Results

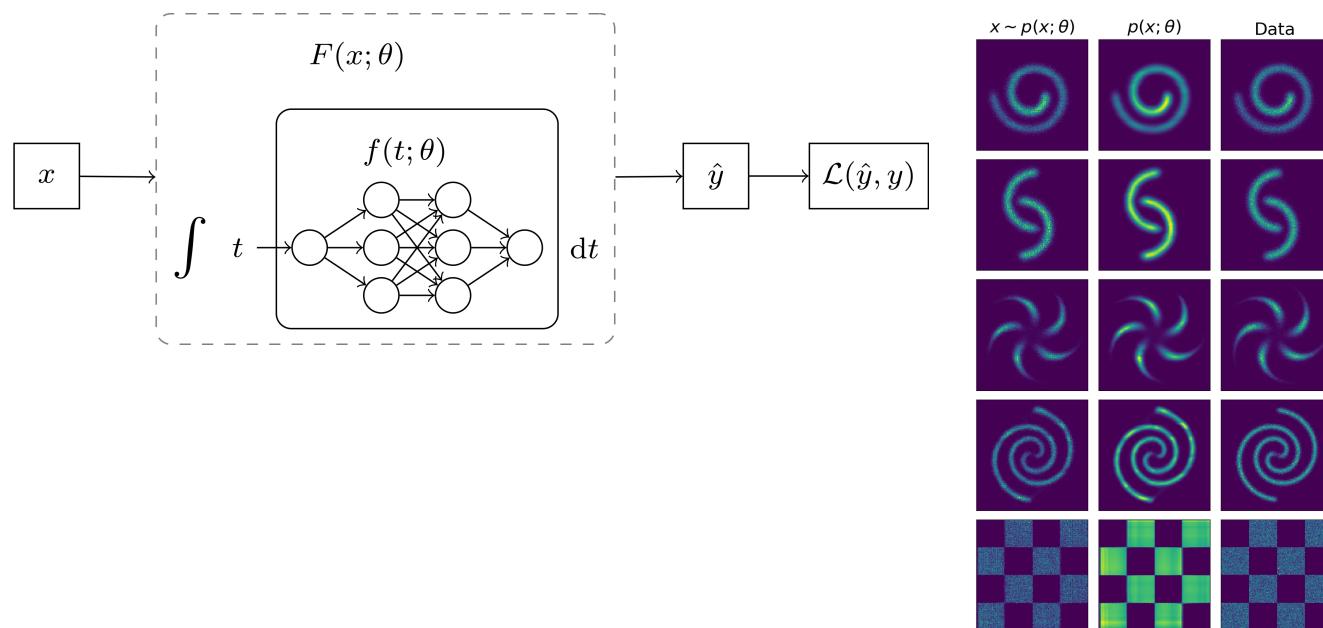
## MNIST



- Learnt by maximum likelihood on the training samples.
- Data dimensionality is 784.
- Samples are generated by sampling noise and inverting the model.
- The model has no knowledge about the structure of the images, i.e.:  
$$p(\mathbf{x}) = p(x_1, \dots, x_{784}) = p(x_1) \prod_{i=2}^{784} p(x_i | x_1, \dots, x_{i-1})$$

# Take Home Messages

- Any monotonic function can be modeled by a neural network that represents the function derivative.
- The Backward pass is memory efficient thanks to the Leibnitz rule.
- UMNN can be used as a building block of autoregressive bijective maps and provide a state of the art density estimator.



*Fin*

# Experimental Results

## Density Estimation Benchmarks

Dataset	POWER	GAS	HEPMASS	MINIBOONE	BSDS300	MNIST
RealNVP - <a href="#">Dinh et al. [2017]</a>	$-0.17 \pm .01$	$-8.33 \pm .14$	$18.71 \pm .02$	$13.55 \pm .49$	$-153.28 \pm 1.78$	-
(a) Glow - <a href="#">Kingma and Dhariwal [2018]</a>	$-0.17 \pm .01$	$-8.15 \pm .40$	$19.92 \pm .08$	$11.35 \pm .07$	$-155.07 \pm .03$	-
FFJORD - <a href="#">Grathwohl et al. [2018]</a>	$\underline{-0.46} \pm .01$	$\underline{-8.59} \pm .12$	$\underline{14.92} \pm .08$	$\underline{10.43} \pm .04$	$\underline{-157.40} \pm .19$	-
MADE - <a href="#">Germain et al. [2015]</a>	$3.08 \pm .03$	$-3.56 \pm .04$	$20.98 \pm .02$	$15.59 \pm .50$	$-148.85 \pm .28$	$2.04 \pm .01$
(b) MAF - <a href="#">Papamakarios et al. [2017]</a>	$-0.24 \pm .01$	$-10.08 \pm .02$	$17.70 \pm .02$	$11.75 \pm .44$	$-155.69 \pm .28$	$1.89 \pm .01$
TAN - <a href="#">Oliva et al. [2018]</a>	$\underline{-0.60} \pm .01$	$\underline{-12.06} \pm .02$	$\underline{13.78} \pm .02$	$\underline{11.01} \pm .48$	$\underline{-159.80} \pm .07$	$\underline{1.19}$
NAF - <a href="#">Huang et al. [2018]</a>	$-0.62 \pm .01$	$-11.96 \pm .33$	$15.09 \pm .40$	$\underline{8.86} \pm .15$	$-157.73 \pm .30$	-
(c) B-NAF - <a href="#">De Cao et al. [2019]</a>	$-0.61 \pm .01$	$\underline{-12.06} \pm .09$	$14.71 \pm .38$	$8.95 \pm .07$	$-157.36 \pm .03$	-
SOS - <a href="#">Jaini et al. [2019]</a>	$-0.60 \pm .01$	$-11.99 \pm .41$	$15.15 \pm .1$	$8.90 \pm .11$	$-157.48 \pm .41$	$1.81$
<b>UMNN-MAF (ours)</b>	$\underline{\textbf{-0.63}} \pm .01$	$-10.89 \pm .7$	$\underline{13.99} \pm .21$	$9.67 \pm .13$	$\underline{-157.98} \pm .01$	$\underline{\textbf{1.13}} \pm .02$

- Learnt by maximum likelihood on the training samples.

# Experimental Results

## Density Estimation Benchmarks

Dataset	POWER	GAS	HEPMASS	MINIBOONE	BSDS300	MNIST
RealNVP - <a href="#">Dinh et al. [2017]</a>	$-0.17 \pm .01$	$-8.33 \pm .14$	$18.71 \pm .02$	$13.55 \pm .49$	$-153.28 \pm 1.78$	-
(a) Glow - <a href="#">Kingma and Dhariwal [2018]</a>	$-0.17 \pm .01$	$-8.15 \pm .40$	$19.92 \pm .08$	$11.35 \pm .07$	$-155.07 \pm .03$	-
FFJORD - <a href="#">Grathwohl et al. [2018]</a>	$\underline{-0.46} \pm .01$	$\underline{-8.59} \pm .12$	$\underline{14.92} \pm .08$	$\underline{10.43} \pm .04$	$\underline{-157.40} \pm .19$	-
MADE - <a href="#">Germain et al. [2015]</a>	$3.08 \pm .03$	$-3.56 \pm .04$	$20.98 \pm .02$	$15.59 \pm .50$	$-148.85 \pm .28$	$2.04 \pm .01$
(b) MAF - <a href="#">Papamakarios et al. [2017]</a>	$-0.24 \pm .01$	$-10.08 \pm .02$	$17.70 \pm .02$	$11.75 \pm .44$	$-155.69 \pm .28$	$1.89 \pm .01$
TAN - <a href="#">Oliva et al. [2018]</a>	$\underline{-0.60} \pm .01$	$\underline{-12.06} \pm .02$	$\underline{13.78} \pm .02$	$\underline{11.01} \pm .48$	$\underline{-159.80} \pm .07$	$1.19$
NAF - <a href="#">Huang et al. [2018]</a>	$-0.62 \pm .01$	$-11.96 \pm .33$	$15.09 \pm .40$	$\underline{8.86} \pm .15$	$-157.73 \pm .30$	-
(c) B-NAF - <a href="#">De Cao et al. [2019]</a>	$-0.61 \pm .01$	$\underline{-12.06} \pm .09$	$14.71 \pm .38$	$8.95 \pm .07$	$-157.36 \pm .03$	-
SOS - <a href="#">Jaini et al. [2019]</a>	$-0.60 \pm .01$	$-11.99 \pm .41$	$15.15 \pm .1$	$8.90 \pm .11$	$-157.48 \pm .41$	$1.81$
<b>UMNN-MAF (ours)</b>	$\underline{\textbf{-0.63}} \pm .01$	$-10.89 \pm .7$	$\underline{13.99} \pm .21$	$9.67 \pm .13$	$\underline{-157.98} \pm .01$	$\textbf{1.13} \pm .02$

- Learnt by maximum likelihood on the training samples.

# Experimental Results

## Density Estimation Benchmarks

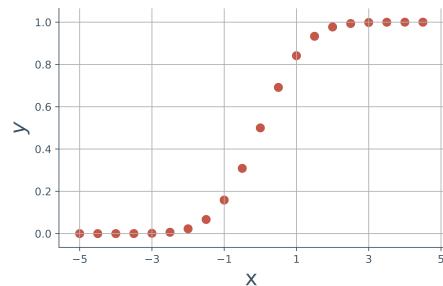
Dataset	POWER	GAS	HEPMASS	MINIBOONE	BSDS300	MNIST
RealNVP - Dinh et al. [2017]	$-0.17 \pm .01$	$-8.33 \pm .14$	$18.71 \pm .02$	$13.55 \pm .49$	$-153.28 \pm 1.78$	-
(a) Glow - Kingma and Dhariwal [2018]	$-0.17 \pm .01$	$-8.15 \pm .40$	$19.92 \pm .08$	$11.35 \pm .07$	$-155.07 \pm .03$	-
FFJORD - Grathwohl et al. [2018]	$\underline{-0.46} \pm .01$	$\underline{-8.59} \pm .12$	$\underline{14.92} \pm .08$	$\underline{10.43} \pm .04$	$\underline{-157.40} \pm .19$	-
MADE - Germain et al. [2015]	$3.08 \pm .03$	$-3.56 \pm .04$	$20.98 \pm .02$	$15.59 \pm .50$	$-148.85 \pm .28$	$2.04 \pm .01$
(b) MAF - Papamakarios et al. [2017]	$-0.24 \pm .01$	$-10.08 \pm .02$	$17.70 \pm .02$	$11.75 \pm .44$	$-155.69 \pm .28$	$1.89 \pm .01$
TAN - Oliva et al. [2018]	$\underline{-0.60} \pm .01$	$\underline{-12.06} \pm .02$	$\underline{13.78} \pm .02$	$\underline{11.01} \pm .48$	$\underline{-159.80} \pm .07$	$\underline{1.19}$
NAF - Huang et al. [2018]	$-0.62 \pm .01$	$-11.96 \pm .33$	$15.09 \pm .40$	$8.86 \pm .15$	$-157.73 \pm .30$	-
(c) B-NAF - De Cao et al. [2019]	$-0.61 \pm .01$	$\underline{-12.06} \pm .09$	$14.71 \pm .38$	$8.95 \pm .07$	$-157.36 \pm .03$	-
SOS - Jaini et al. [2019]	$-0.60 \pm .01$	$-11.99 \pm .41$	$15.15 \pm .1$	$8.90 \pm .11$	$-157.48 \pm .41$	$1.81$
<b>UMNN-MAF (ours)</b>	$\underline{\textbf{-0.63}} \pm .01$	$-10.89 \pm .7$	$\underline{13.99} \pm .21$	$9.67 \pm .13$	$\underline{-157.98} \pm .01$	$\underline{\textbf{1.13}} \pm .02$

- Learnt by maximum likelihood on the training samples.

# UMNN to the rescue 😱

In Practice:

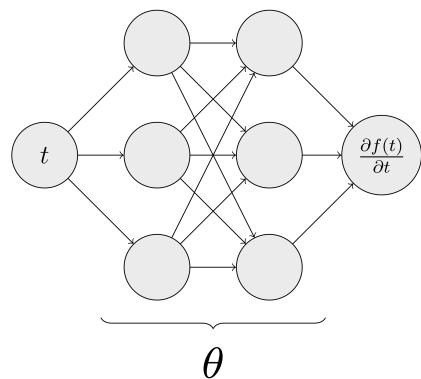
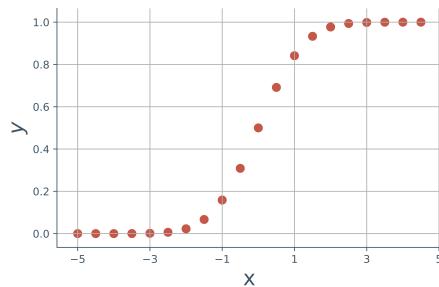
A possible solution: to numerically approximate the derivative at each point.



# UMNN to the rescue 😱

In Practice:

A possible solution: to numerically approximate the derivative at each point.



A better solution: to integrate the model of the derivative.  $\hat{y} = \int_0^x \frac{\partial f(t)}{\partial t} dt + f(0)$