

Normalizing Flows

Didrik Nielsen



Where are flows useful?

Anywhere you need a flexible density $p(\mathbf{x})$ with

Sampling:

$$\mathbf{x} \sim p(\mathbf{x})$$

Density:

$$\mathbf{x} \mapsto \log p(\mathbf{x})$$



Why Flows?

	Flow	ARM	VAE	GAN	EBM	DDPM
Density	Exact Fast	Exact Fast	Approx. Fast	NA	Unnorm.	Approx. Slow
Sampling	Fast	Slow	Fast	Fast	Slow & Approx.	Slow

History

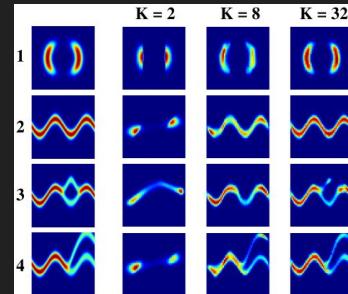
Normalizing Flows: Originates from (Tabak & Turner 2013; Tabak & Vanden-Eijnden 2010)
Describes the **change of probability density** through a series of **invertible maps**

Popularized by (Rezende & Mohammed 2015):

Variational Inference with Normalizing Flows

Danilo Jimenez Rezende
Shakir Mohamed
Google DeepMind, London

DANILOR@GOOGLE.COM
SHAKIR@GOOGLE.COM



Model	$-\ln p(\mathbf{x})$
DLGM diagonal covariance	≤ 89.9
DLGM+NF (k = 10)	≤ 87.5
DLGM+NF (k = 20)	≤ 86.5
DLGM+NF (k = 40)	≤ 85.7
DLGM+NF (k = 80)	≤ 85.1

First deep generative flows: NICE (Dinh et al. 2014); RealNVP (Dinh et al. 2016):

NICE: NON-LINEAR INDEPENDENT COMPONENTS ESTIMATION

Laurent Dinh David Krueger Yoshua Bengio*
Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

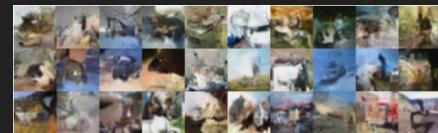


DENSITY ESTIMATION USING REAL NVP

Laurent Dinh*
Montreal Institute for Learning Algorithms
University of Montreal
Montreal, QC H3T1J4

Jascha Sohl-Dickstein
Google Brain

Samy Bengio
Google Brain



Outline

1. Framework
2. Coupling Flows,
Autoregressive Flows
3. Continuous-time Flows,
Residual Flows
4. Discrete Flows,
Surjective Flows

The Flow Framework

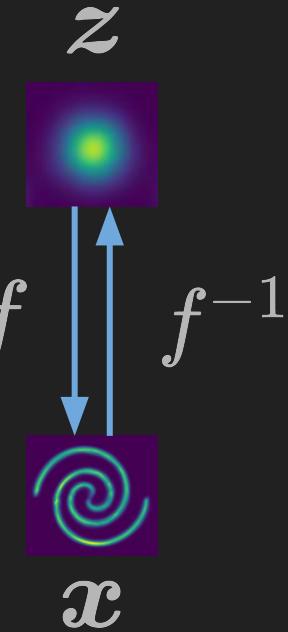
Construct $p(\mathbf{x})$ using

- Base distribution $p(\mathbf{z})$
- Bijective mapping f

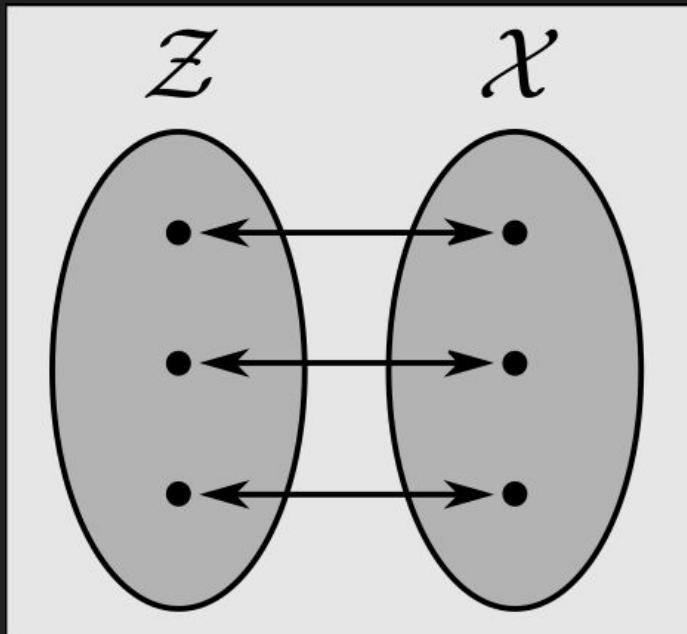
Sampling:

$$\mathbf{z} \sim p(\mathbf{z})$$

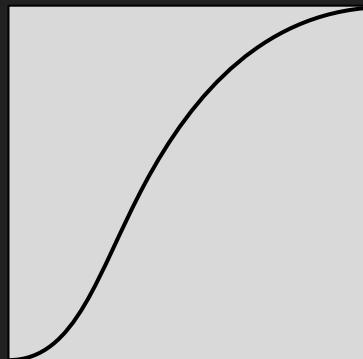
$$\mathbf{x} = f(\mathbf{z})$$



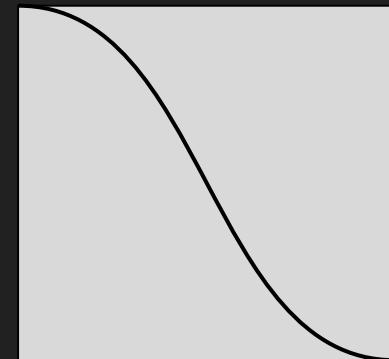
Flows: Invertibility



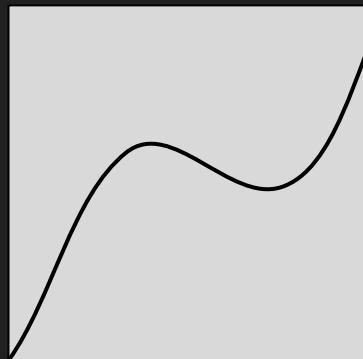
In 1D:
Strictly increasing or decreasing



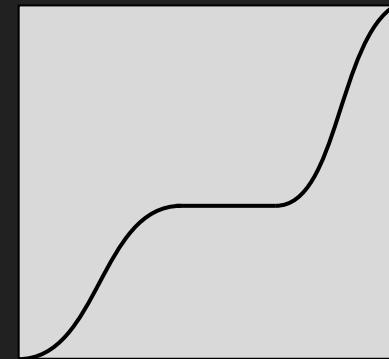
Invertible



Invertible



Not Invertible



Not Invertible

The Flow Framework

Construct $p(\mathbf{x})$ using

- Base distribution $p(\mathbf{z})$
- Bijective mapping f

Sampling:

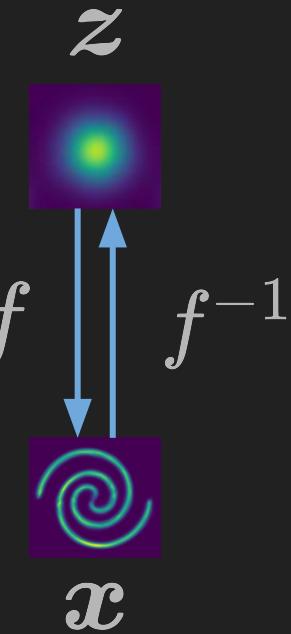
$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = f(\mathbf{z})$$

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

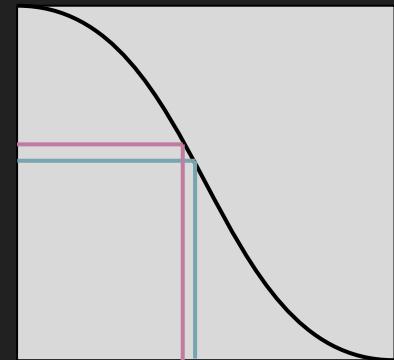
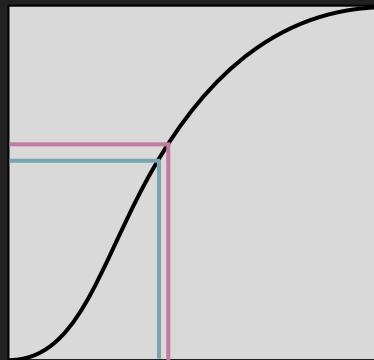
$$\mathbf{z} = f^{-1}(\mathbf{x})$$



Change-of-Variables in 1D

$$\begin{aligned}\int_{x_1}^{x_2} \underline{p(x)dx} &= \int_{z_1}^{z_2} p(z)dz \\ &= \int_{?}^? \underline{p(z) \frac{dz}{dx} dx}\end{aligned}$$

$$\frac{dz}{dx} > 0 \qquad \qquad \frac{dz}{dx} < 0$$



Let's derive it!

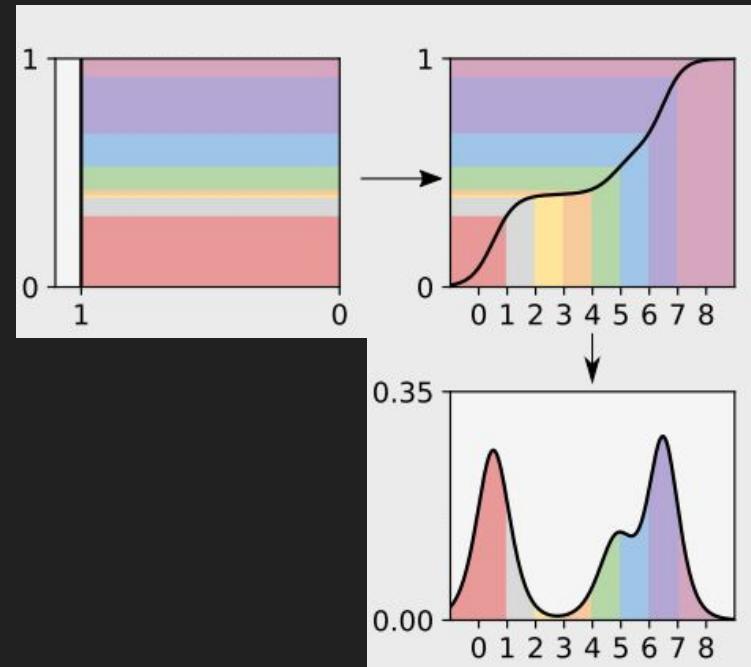
$$p(x) = p(z) \left| \frac{dz}{dx} \right|$$

$$\begin{aligned}\int_{x_1}^{x_2} p(x)dx &= \int_{z_1}^{z_2} p(z)dz \\ &= \int_{x_1}^{x_2} p(z) \frac{dz}{dx} dx \\ &= \int_{x_1}^{x_2} p(z) \left[-\frac{dz}{dx} \right] dx\end{aligned}$$

Flows: 1D Example

Inverse transform sampling:

$$\begin{aligned} p(x) &= \text{Unif}(z|0, 1) \left| \frac{dCDF(x)}{dx} \right| \\ &= \left| \frac{dCDF(x)}{dx} \right| \end{aligned}$$

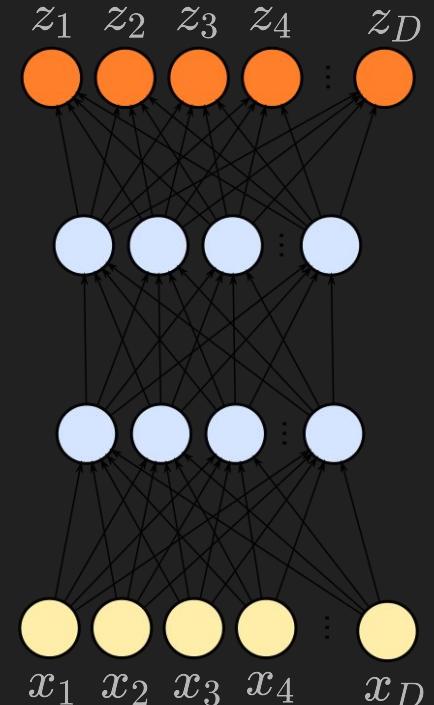


$$p(x) = p(z) \left| \frac{dz}{dx} \right|$$

Change-of-Variables

$$\mathbf{J} = \boxed{\frac{\partial \mathbf{z}}{\partial \mathbf{x}}} = \underbrace{\begin{bmatrix} \frac{\partial z_1}{\partial x_1} & \cdots & \frac{\partial z_1}{\partial x_D} \\ \vdots & & \vdots \\ \frac{\partial z_D}{\partial x_1} & \cdots & \frac{\partial z_D}{\partial x_D} \end{bmatrix}}_{D \times D}$$

$$p(\mathbf{x}) = p(\mathbf{z}) \left| \det \boxed{\frac{\partial \mathbf{z}}{\partial \mathbf{x}}} \right|$$



The Flow Framework

Construct $p(\mathbf{x})$ using

- Base distribution $p(\mathbf{z})$
- Bijective mapping f

Sampling:

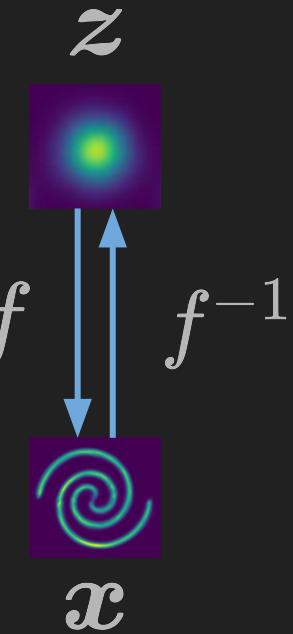
$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = f(\mathbf{z})$$

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

$$\mathbf{z} = f^{-1}(\mathbf{x})$$



The Flow Framework

Sampling:

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\underline{\mathbf{x} = f(\mathbf{z})}$$

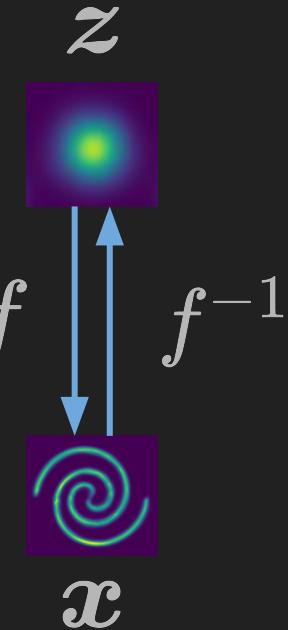
Forward

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

$$\underline{\mathbf{z} = f^{-1}(\mathbf{x})}$$

Inverse



The Flow Framework

Sampling:

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = f_T \circ \dots \circ f_1(\mathbf{z})$$

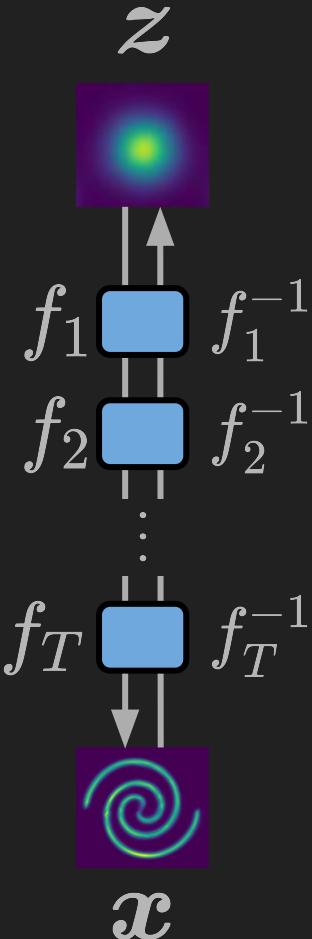
Forward

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \sum_{t=1}^T \log \left| \det \frac{\partial \mathbf{z}_{t-1}}{\partial \mathbf{z}_t} \right|, \quad \text{for } \mathbf{z}_T = \mathbf{x}$$

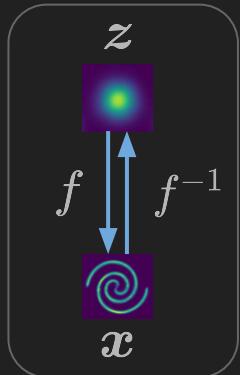
$$\mathbf{z} = f_1^{-1} \circ \dots \circ f_T^{-1}(\mathbf{x})$$

Inverse



Outline

1. Framework



Sampling:

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = \underline{f(\mathbf{z})}$$

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

$$\mathbf{z} = \underline{f^{-1}(\mathbf{x})}$$

2. Coupling Flows, Autoregressive Flows

3. Continuous-time Flows, Residual Flows

4. Discrete Flows, Surjective Flows

How to Build Efficient Flows?

All about developing layers that:

- 1) Are **expressive**
- 2) Are **invertible**
- 3) Have **cheap-to-compute** Jacobian determinants

Det. Identities:

Planar, radial, Sylvester, etc.

Autoregressive:

MAF, IAF, NAF, TAN, etc.

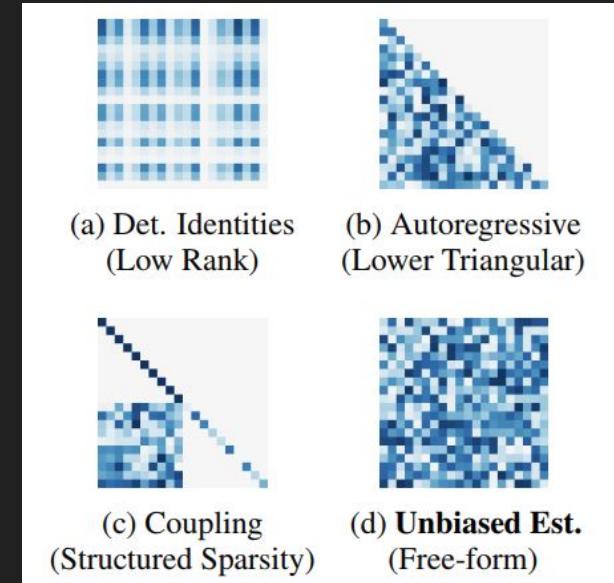
Coupling:

NICE, RealNVP, Glow, Flow++, etc.

Popular!

Unbiased:

FFJORD, Residual Flows, etc.



Taken from (Chen et al., 2019)

Coupling Flows

Forward:

$$\mathbf{x}_{1:d} = \mathbf{z}_{1:d}$$

$$\mathbf{x}_{d+1:D} = (\mathbf{z}_{d+1:D} - \mu_{d+1:D}) \cdot e^{-\alpha_{d+1:D}}$$

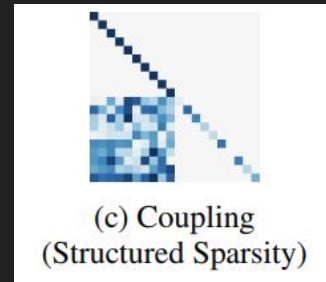
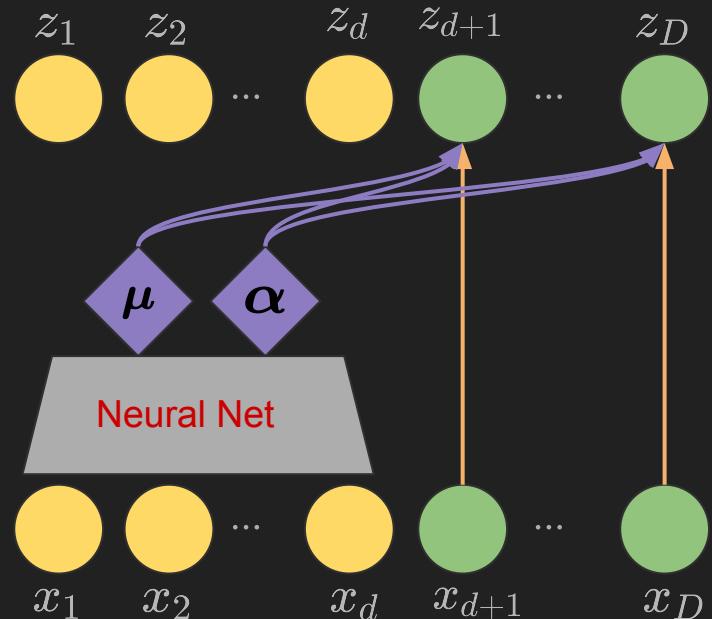
Inverse:

$$\mathbf{z}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{z}_{d+1:D} = \mathbf{x}_{d+1:D} \cdot e^{\alpha_{d+1:D}} + \mu_{d+1:D}$$

Jac. Det:

$$\log |\det J| = \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right| = \sum_{i=d+1}^D \alpha_i$$



(c) Coupling
(Structured Sparsity)

Are Coupling Layers Enough?

Problem?

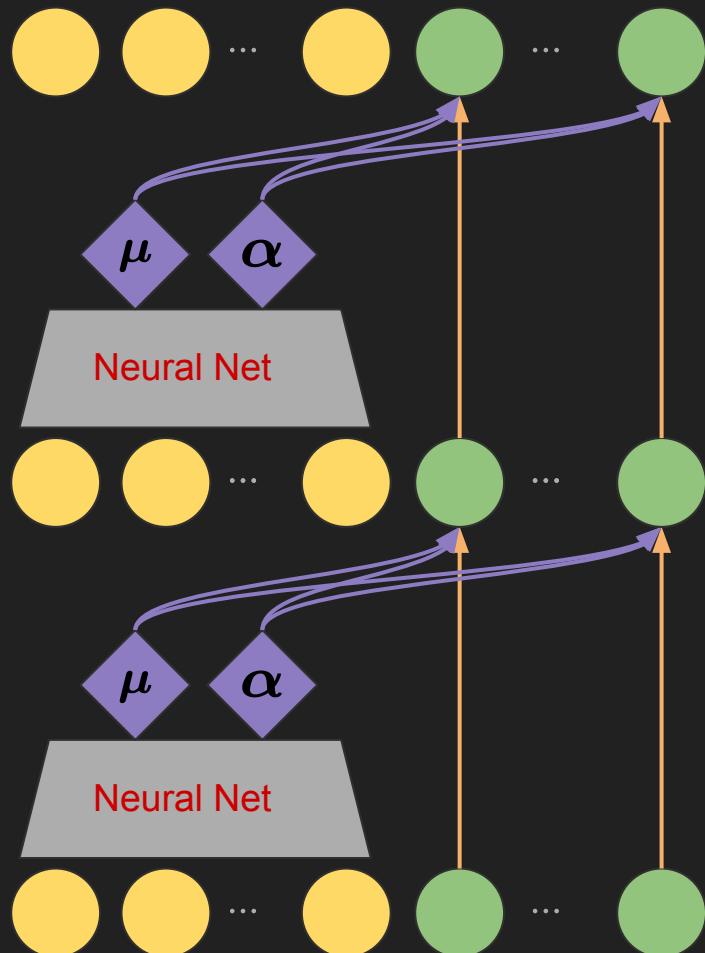
No mixing, need permutation layers!

Reverse:

$$x = \begin{bmatrix} 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \end{bmatrix} z$$

Shuffle:

$$x = \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} z$$



Time to Code!

Tasks

We will implement a coupling flow like RealNVP/Glow and fit it to 2D toy data.

Todo:

1. Implement the `Coupling` bijection class
2. Implement the `make_net` function
3. Train flow and sample
4. **Bonus:** Study the effect of some parameter (e.g. flow layers, NN layers, NN dim, NN activation, etc.) and improve performance.

Useful functions:

Problem 1: `torch.chunk`, `torch.cat`

Problem 2: `torch.nn.Sequential`

Forward:

$$\mathbf{x}_{1:d} = \mathbf{z}_{1:d}$$

$$\mathbf{x}_{d+1:D} = (\mathbf{z}_{d+1:D} - \mu_{d+1:D}) \cdot e^{-\alpha_{d+1:D}}$$

Inverse:

$$\mathbf{z}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{z}_{d+1:D} = \mathbf{x}_{d+1:D} \cdot e^{\alpha_{d+1:D}} + \mu_{d+1:D}$$

Jac. Det:

$$\log |\det J| = \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right| = \sum_{i=d+1}^D \alpha_i$$

`.inverse(self, z)`

⋮

return \mathbf{x}

(B, 2)

`.forward(self, x)`

⋮

return \mathbf{z}, ldj

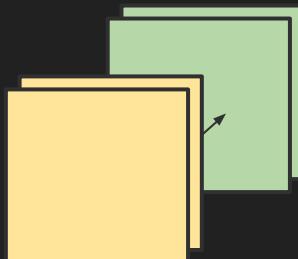
(B, 2) (B,)

(Taken from Kingma & Dhariwal 2018)

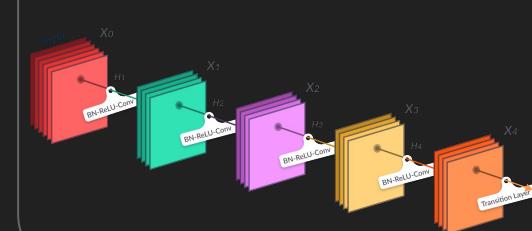
Coupling Flows on Images

While toy data had shape (B,2),
e.g. CIFAR-10 has shape (B,3,32,32)

Split along channels:



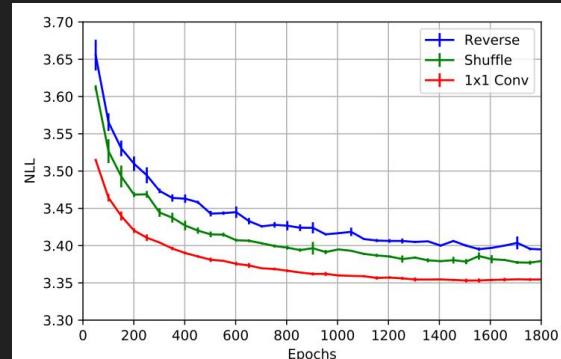
Use e.g. CNNs:



Need dequantization
since images are
discrete in $\{0, 1, \dots, 255\}$



Permutation using Conv1x1:



How to Build Efficient Flows?

All about developing layers that:

- 1) Are **expressive**
- 2) Are **invertible**
- 3) Have **cheap-to-compute** Jacobian determinants

Det. Identities:

Planar, radial, Sylvester, etc.

Autoregressive:

MAF, IAF, NAF, TAN, etc.

Related to
AR models

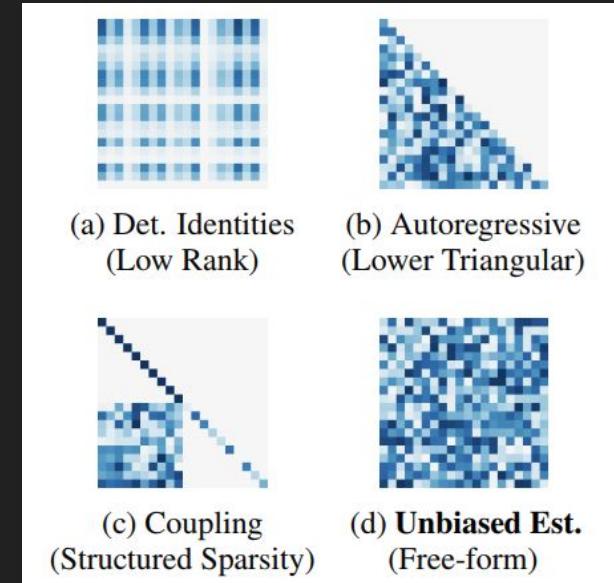
Popular!

Coupling:

NICE, RealNVP, Glow, Flow++, etc.

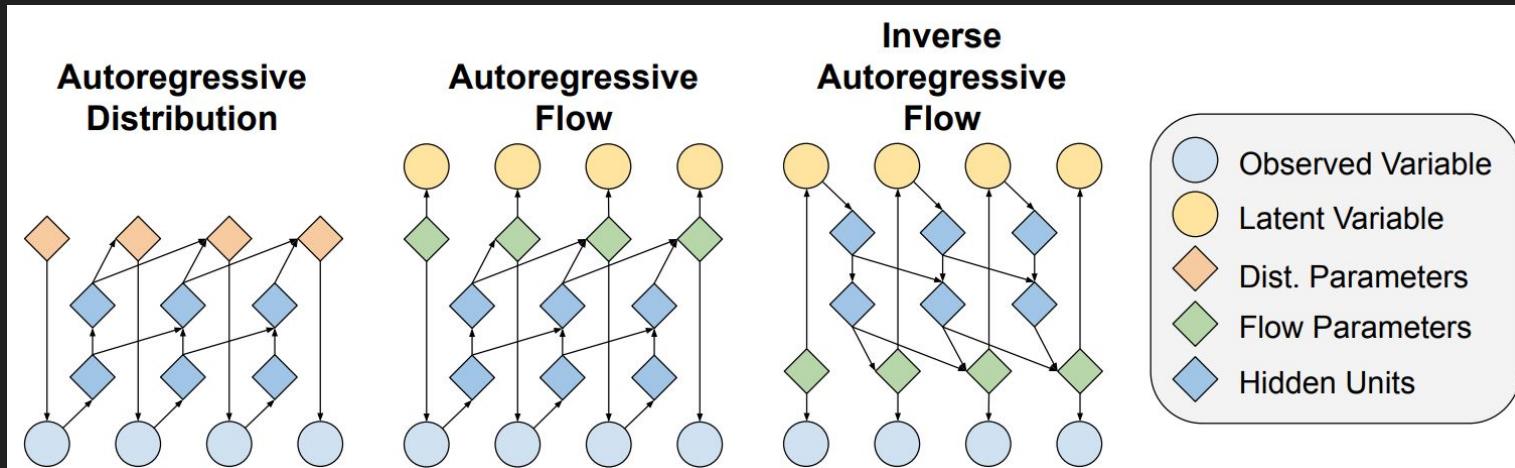
Unbiased:

FFJORD, Residual Flows, etc.



Taken from (Chen et al., 2019)

Autoregressive Flows



Equivalent for continuous distributions
(Papamakarios et al. 2017)

Sampling: D passes (slow)
Density: 1 pass (fast)

Inverse possible for flows
(Kingma et al. 2016)

Sampling: 1 pass (fast)
Density: D passes (slow)

(for comparison)
Coupling flows:

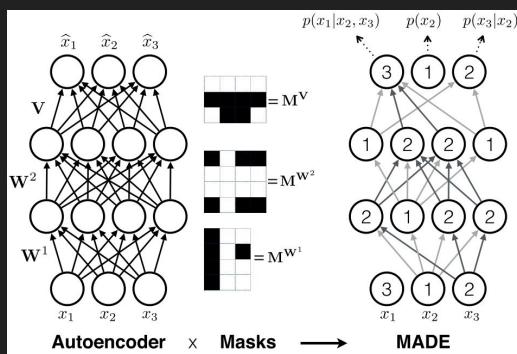
Sampling: 1 pass (fast)
Density: 1 pass (fast)

RNNs (e.g. LSTM, GRU)
require sequential computation

Masked Autoregressive NNs

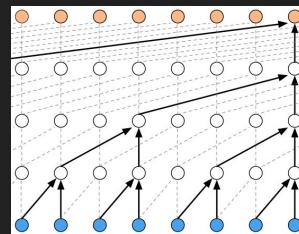
Vectors (B,D)

MADE
(Larochelle et al. 2015)



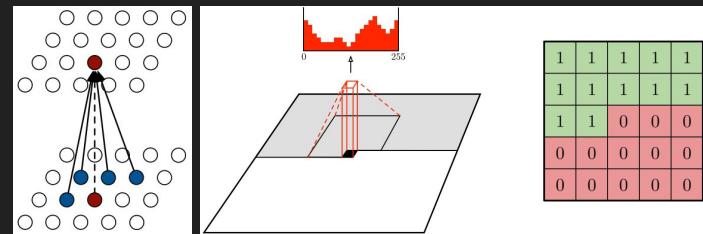
Sequences (B,D,L)

WaveNet
(van den Oord et al. 2016)



Images (B,C,H,W)

PixelCNN
(van den Oord et al. 2016)



...

Transformer
(Vaswani et al. 2017)

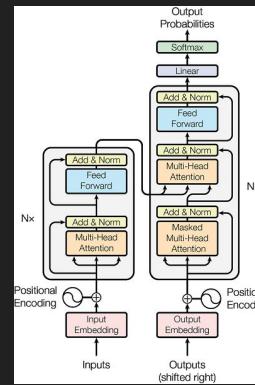
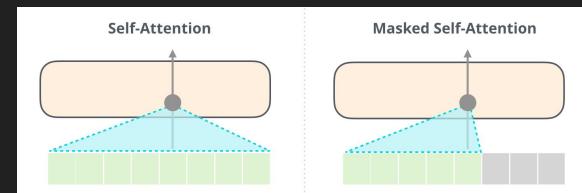
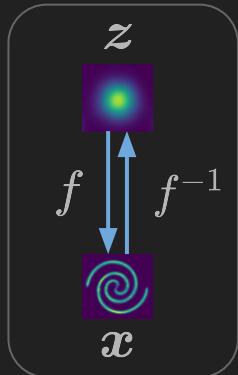


Image Transformer
(Parmar et al. 2018)



Outline

1. Framework



Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$
$$\mathbf{z} = \underline{f^{-1}(\mathbf{x})}$$

2. Coupling Flows, Autoregressive Flows



Sampling: 1 pass (fast)
Density: 1 pass (fast)



Sampling: D passes (slow)
Density: 1 pass (fast)

3. Continuous-time Flows, Residual Flows

4. Discrete Flows, Surjective Flows

Continuous-time Flows

Neural ODE (Chen et al. 2018):

$$\frac{\partial \mathbf{z}(t)}{\partial t} = f_{\theta}(\mathbf{z}(t), t)$$

Instantaneous Change-of-Variables formula:

$$\frac{\partial \log p(\mathbf{z}(t))}{\partial t} = -\text{Tr} \left(\frac{\partial f}{\partial \mathbf{z}(t)} \right)$$

FFJORD (Grathwohl et al. 2018):

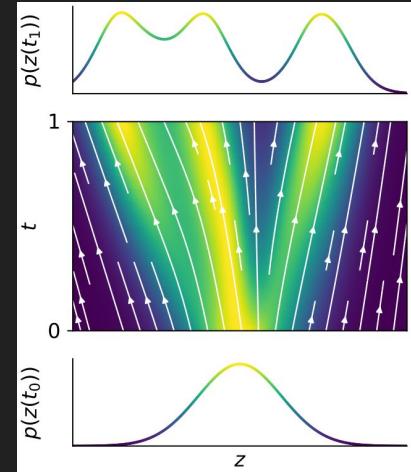
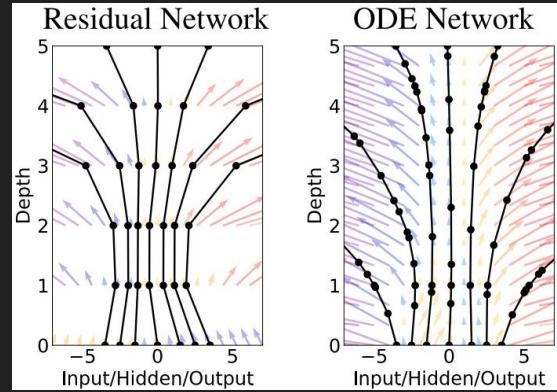
Use Hutchinson's Trace Estimator:

$$\text{Tr} \left(\frac{\partial f}{\partial \mathbf{z}(t)} \right) = \mathbb{E}_{p(\boldsymbol{\epsilon})} \left[\boldsymbol{\epsilon}^T \frac{\partial f}{\partial \mathbf{z}(t)} \boldsymbol{\epsilon} \right] \quad \text{for} \quad \mathbb{E} [\boldsymbol{\epsilon}] = \mathbf{0}, \text{Cov} [\boldsymbol{\epsilon}] = \mathbf{I}$$

Jacobian $\frac{\partial f}{\partial \mathbf{z}(t)}$ costs $O(D^2)$

vjp $\boldsymbol{\epsilon}^T \frac{\partial f}{\partial \mathbf{z}(t)}$ costs $O(D)$

Connects to score-based / diffusion models!
(Song et al. 2020)



Residual Flows

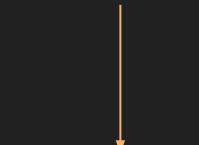
i-ResNets (Behrmann et al. 2019):

$$\mathbf{z} = f(\mathbf{x}) = \mathbf{x} + h(\mathbf{x}) \quad \text{for contractive } h$$

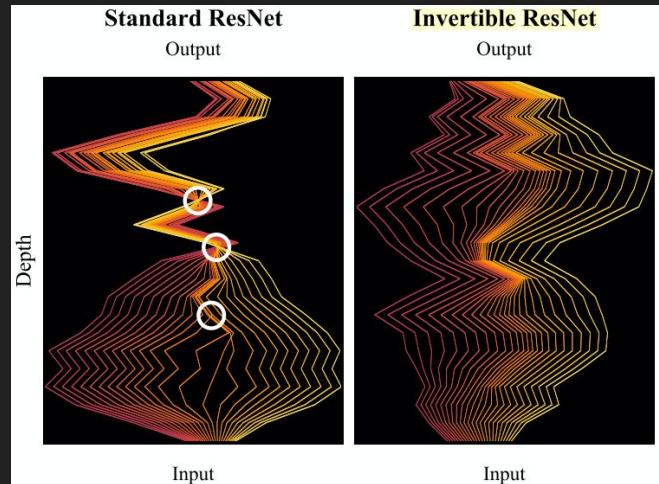
$$\log J_f = \text{Tr} \left(\sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \mathbf{J}_h^k \right)$$

Hutchinson's
Trace Estimator

Truncate
(introduces bias)



Enforce using
SpectralNorm

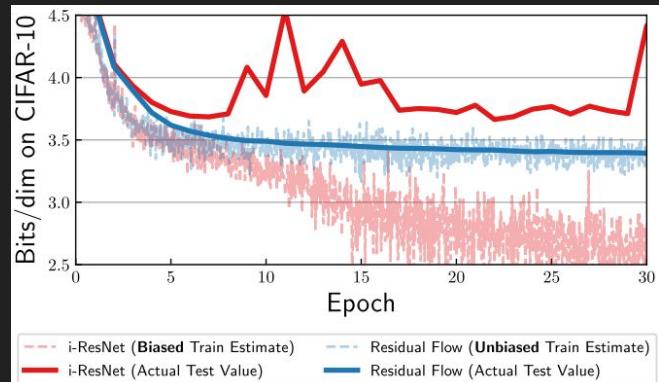


Residual Flows (Chen et al. 2019):

Replace truncation by (unbiased) Russian Roulette Estimator:

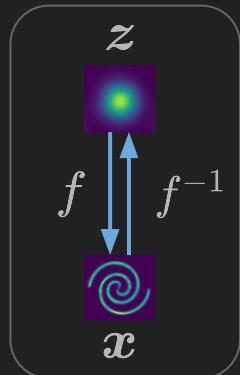
$$\sum_{k=1}^{\infty} \Delta_k = \mathbb{E}_{p(N)} \left[\sum_{k=1}^N \frac{\Delta_k}{P(N \geq k)} \right] \quad \text{for } p(N) \text{ w/ integer support}$$

+ introduce memory-efficient gradient estimators



Outline

1. Framework



Sampling:

$$\mathbf{z} \sim p(\mathbf{z})$$

$$\mathbf{x} = \underline{f(\mathbf{z})}$$

Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$

$$\mathbf{z} = \underline{f^{-1}(\mathbf{x})}$$

2. Coupling Flows, Autoregressive Flows



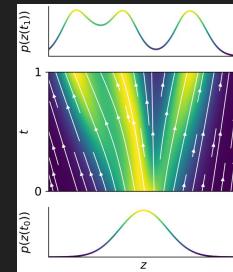
Sampling: 1 pass (**fast**)
Density: 1 pass (**fast**)



Sampling: D passes (**slow**)
Density: 1 pass (**fast**)

Sampling: 1 pass (**fast**)
Density: D passes (**slow**)

3. Continuous-time Flows, Residual Flows



ODE



Invertible
ResNet

4. Discrete Flows, Surjective Flows

Discrete Flows

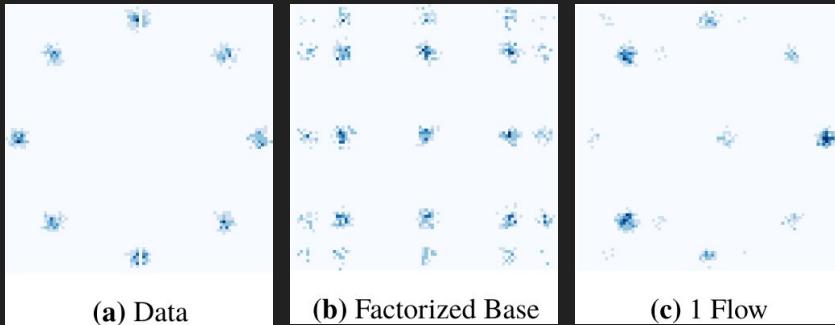
No volume: $\log p(\mathbf{x}) = \log p(\mathbf{z})$

Discrete Flows (Tran et al. 2019):

$$z_d = (\mu_d + \sigma_d x_d) \bmod K$$

Use Straight-Through Estimator:

Forward: $\mu_d = \text{one_hot}(\text{argmax}(\theta_d))$
Backward: $\text{softmax}(\theta_d/\tau)$

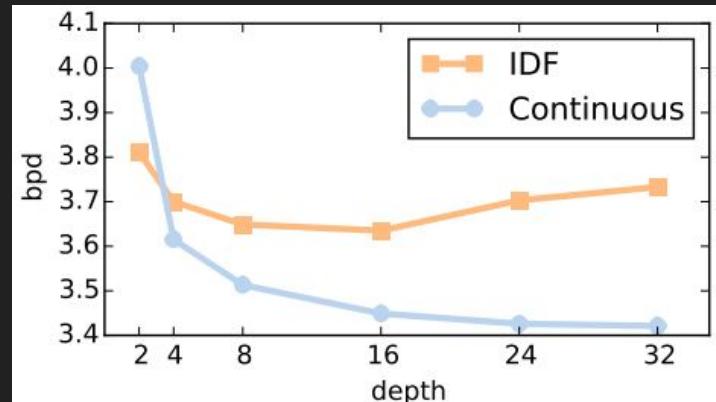


Integer Discrete Flows (Hoogeboom et al. 2019):

$$z_d = x_d + \mu_d$$

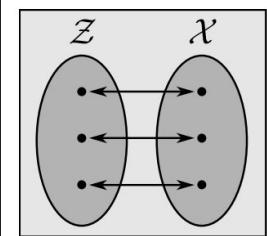
Also use Straight-Through Estimator:

Forward: $\mu_d = \lfloor \theta_d \rfloor$
Backward: θ_d

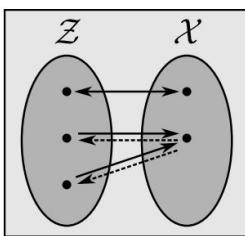


Surjective Flows

SurVAE Flows (Nielsen et al. 2020):



Bijective



Surjective



Stochastic

Examples:

$x = \text{round}(z)$

Dequantization (Uria et al., 2013; Ho et al. 2019)

$x = z[:n]$

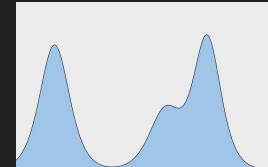
Augmented Normalizing Flows (Huang et al., 2020), VFlow (Chen et al., 2020)

$x = \text{argmax}(z)$

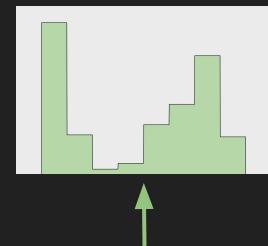
Argmax Flows (Hoogeboom et al., 2021)

⋮

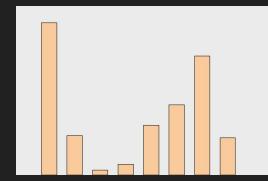
Flow



Dequantized Data



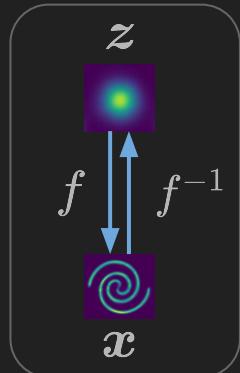
Data



heartedness frege thematically inferred by the famous existence of a function f from the laplace definition we can analyze a definition of binary operations with additional size so their functionality cannot be viewed here there is no change because its total cost of learning objects from language to platonic linguistics examines why animate to indicate wild amphibious substances animal and marine life constituents of animals and bird sciences medieval biology biology and central medicine full discovery re

Outline

1. Framework



Density:

$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}} \right|$$
$$\mathbf{z} = \underline{f^{-1}(\mathbf{x})}$$

2. Coupling Flows, Autoregressive Flows



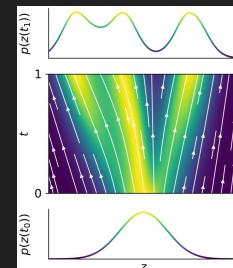
Sampling: 1 pass (**fast**)
Density: 1 pass (**fast**)



Sampling: D passes (**slow**)
Density: 1 pass (**fast**)

Sampling: 1 pass (**fast**)
Density: D passes (**slow**)

3. Continuous-time Flows, Residual Flows



ODE



Invertible
ResNet

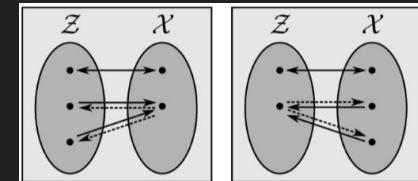
4. Discrete Flows, Surjective Flows

Discrete Flows:

$$\log p(\mathbf{x}) = \log p(\mathbf{z})$$
$$\mathbf{z} = f^{-1}(\mathbf{x})$$

Straight-Through Estimator

Surjective Flows:



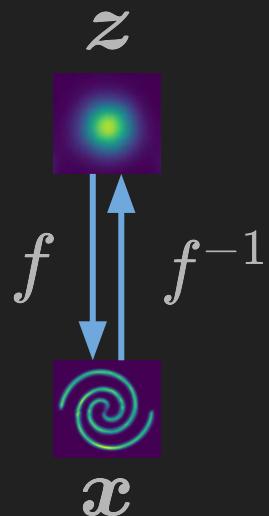
$$\log p(\mathbf{x}) = \log p(\mathbf{z}) + \mathcal{V}(\mathbf{x}, \mathbf{z})$$

Bonus: Variational Inference with Flows

Generative Modeling

Train to learn
data distribution $q(\mathbf{x})$

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

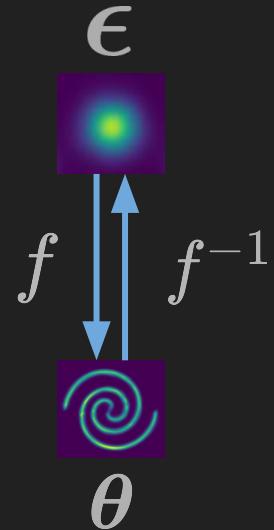


Minimize:
 $\mathbb{D}_{KL} [q(\mathbf{x}) || p(\mathbf{x})]$

Variational Inference

Train to approximate
posterior distribution $p(\boldsymbol{\theta} | \mathcal{D})$

$$q(\boldsymbol{\theta})$$



Minimize:
 $\mathbb{D}_{KL} [q(\boldsymbol{\theta}) || p(\boldsymbol{\theta} | \mathcal{D})]$

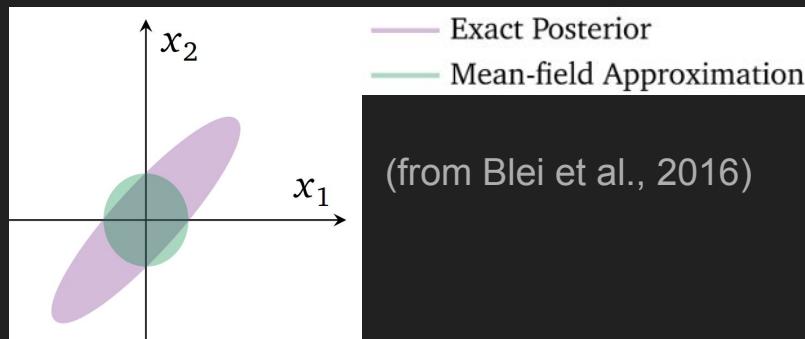
Variational Inference

Learn posterior approximation:

$$q_{\lambda}(\boldsymbol{\theta}) \approx p(\boldsymbol{\theta}|\mathcal{D})$$

By minimizing the KL:

$$\mathbb{D}_{KL} [q_{\lambda}(\boldsymbol{\theta}) \| p(\boldsymbol{\theta}|\mathcal{D})]$$



Example: Mean-field approximation:

Factorized distribution:

$$q_{\lambda}(\boldsymbol{\theta}) = \prod_{d=1}^D q_{\lambda_d}(\theta_d)$$

Common choice:

$$q_{\lambda}(\boldsymbol{\theta}) = \prod_{d=1}^D \mathcal{N}(\theta_d | \mu_d, \sigma_d^2)$$

$$\boldsymbol{\lambda} = \{\boldsymbol{\mu}, \boldsymbol{\sigma}^2\}$$

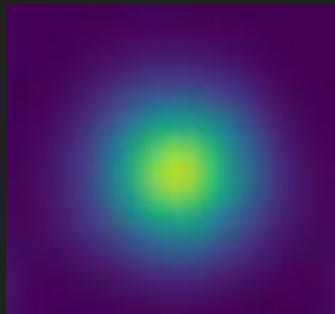
Variational Inference with Flows

Mean-field Gaussian:

$$\epsilon \sim \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})$$

$$\theta = \mu + \sigma \odot \epsilon$$

$$q(\theta) = \mathcal{N}(\theta | \mu, \text{diag}(\sigma^2))$$



Flow:

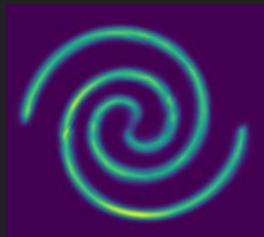
$$\epsilon \sim \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})$$

$$\theta = f_{\lambda}(\epsilon)$$

$$q(\theta) = q(\epsilon) |\det J|$$



Which Flows to use for VI?



$$\epsilon \sim \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})$$
$$\theta = f_{\lambda}(\epsilon)$$

$$q(\theta) = q(\epsilon) |\det \mathbf{J}|$$

Desired properties:

Fast sampling

Density for generated samples

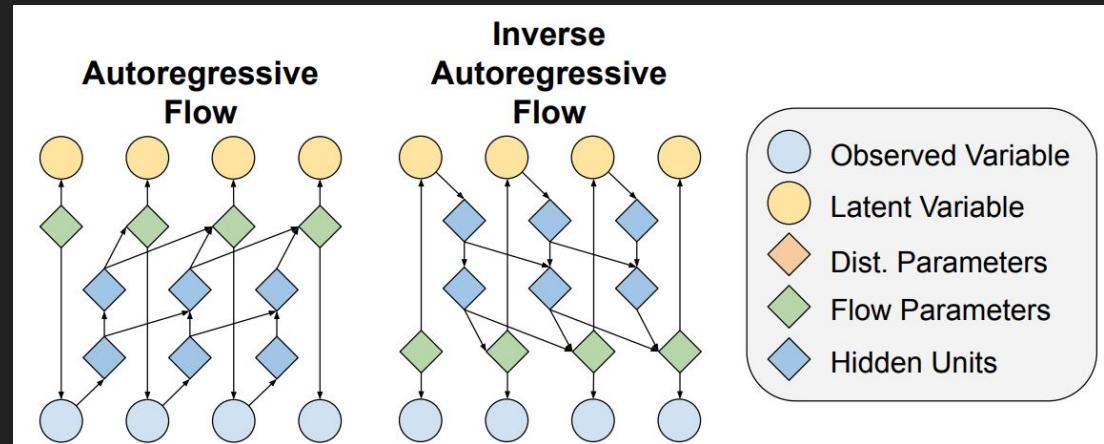
Autoregressive Flow

Forward: D passes (slow)
Inverse: 1 pass (fast)

Inverse Autoregressive Flow

Forward: 1 pass (fast)
Inverse: D passes (slow)

We only need the forward direction!



Other Flows for VI?

Variational Inference with Normalizing Flows

Danilo Jimenez Rezende
Shakir Mohamed
Google DeepMind, London

DANILOR@GOOGLE.COM
SHAKIR@GOOGLE.COM

Det. Identities:

Planar, radial **Sylvester**, etc.

Autoregressive:

MAF, **IAF**, NAF, TAN, etc.

Coupling:

NICE, RealNVP, Glow, Flow++, etc.

Unbiased:

FFJORD, Residual Flows, etc.

Improved Variational Inference with Inverse Autoregressive Flow

Diederik P. Kingma Tim Salimans Rafal Jozefowicz Xi Chen
dpkingma@openai.com tim@openai.com rafal@openai.com peter@openai.com

Ilya Sutskever Max Welling*
ilya@openai.com M.Welling@uva.nl

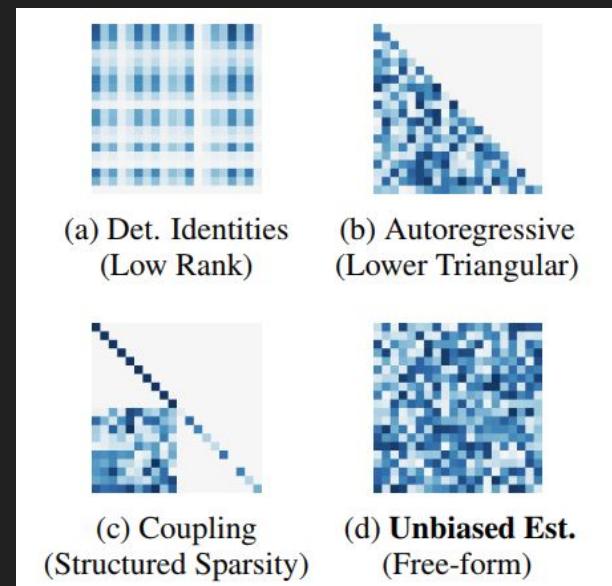
Sylvester Normalizing Flows for Variational Inference

Rianne van den Berg*
University of Amsterdam

Leonard Hasenclever*
University of Oxford

Jakub M. Tomczak
University of Amsterdam

Max Welling
University of Amsterdam



Taken from (Chen et al., 2019)

\mathbf{z}

Example: Variational Autoencoders

Prior: $p(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$

Likelihood: $p(\mathbf{x} | \mathbf{z}) = p(\mathbf{x} | \text{NN}_{\theta}(\mathbf{z}))$

Mean-field:

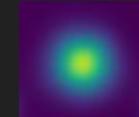
$$q(\mathbf{z} | \mathbf{x}) = \mathcal{N}(\mathbf{z} | \boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma^2}))$$

$$\boldsymbol{\mu}, \boldsymbol{\sigma^2} \leftarrow \text{NN}_{\phi}(\mathbf{x})$$

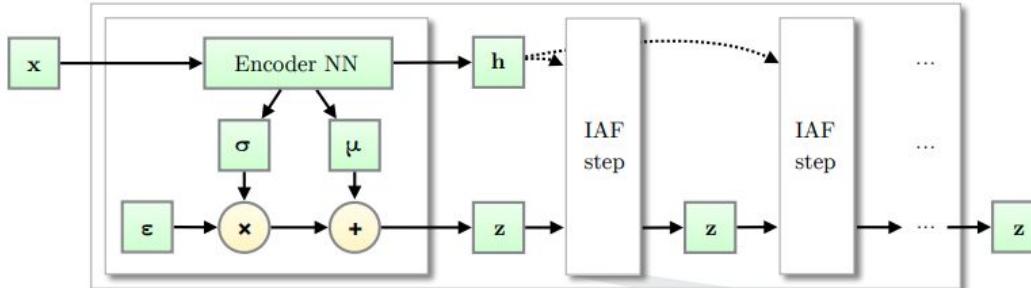
Flow:

$$q(\mathbf{z} | \mathbf{x}) = q(\boldsymbol{\epsilon}) \left| \det \frac{\partial f_{\lambda}(\boldsymbol{\epsilon}; \mathbf{x})}{\partial \boldsymbol{\epsilon}} \right|^{-1}$$

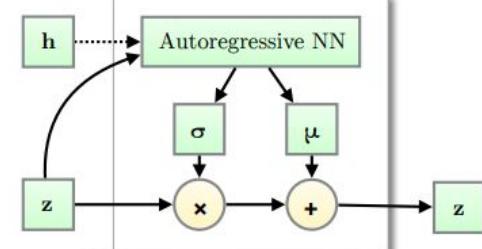
$$p(\mathbf{x} | \mathbf{z}) \quad q(\mathbf{z} | \mathbf{x})$$

 \mathbf{x}

Approximate Posterior with Inverse Autoregressive Flow (IAF)



IAF Step



Example: Bayesian Neural Networks

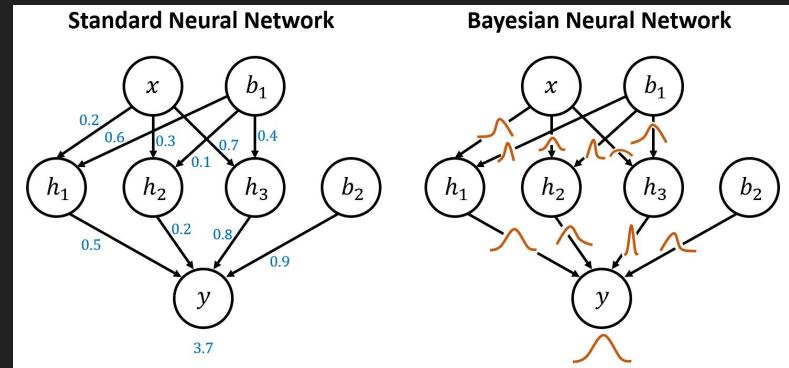
Likelihood: $p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \boldsymbol{\theta})$

$$\mu_n = \text{NN}_{\boldsymbol{\theta}}(\mathbf{x}_n)$$

$$\begin{aligned} &\text{Bern}(y_n|\mu_n) \\ &\mathcal{N}(y_n|\mu_n, \sigma^2) \end{aligned}$$

$$\boldsymbol{\theta} = \{\mathbf{w}_1, \mathbf{b}_1, \dots, \mathbf{w}_L, \mathbf{b}_L\}$$

$$\text{Prior: } p(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, \mathbf{I})$$



Mean-field:

$$q(\boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$$

Flow:

$$q(\boldsymbol{\theta}) = q(\boldsymbol{\epsilon}) \left| \det \frac{\partial f_{\boldsymbol{\lambda}}(\boldsymbol{\epsilon})}{\partial \boldsymbol{\epsilon}} \right|^{-1}$$

Scalable Bayesian Neural Networks

Local Reparameterization Trick (Kingma et al. 2015):

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b} \quad \mathbf{W}, \mathbf{b} \sim q(\mathbf{W}, \mathbf{b})$$

$$q(\mathbf{W}, \mathbf{b}) \sim \mathcal{N} \implies q(\mathbf{h}) \sim \mathcal{N}$$

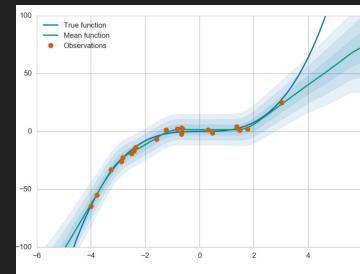
Much lower
gradient variance!

Multiplicative Normalizing Flows (Louizos & Welling 2017):

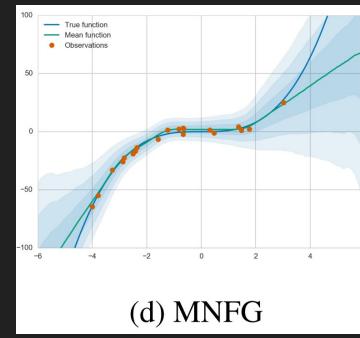
$$q(\mathbf{z}) = q(\boldsymbol{\epsilon}) |\det \mathbf{J}|$$

$$q(\mathbf{W}|\mathbf{z}) = \prod_i \prod_o \mathcal{N}(z_i \mu_{io}, \sigma_{io}^2)$$

- Can still apply LRT
- More flexible q



(c) FFLU



(d) MNFG

Thanks for listening!