# **Deep Generative Models**

## 7. Normalizing Flow Models



• 국가수리과학연구소 산업수학혁신센터 김민중

### Latent Variable Models: Assumption

- Observable variables  $x \in \mathbb{R}^d$
- Latent variables  $z \in \mathbb{R}^h$  (unobservable)

$$p_{data}(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z})$$
  
or 
$$= \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z}$$

### **Example: Mixture of Gaussians**

- Mixture of Gaussians. Bayes net:  $z \rightarrow x$ 
  - $z = Categorical(z|\gamma_1, \cdots, \gamma_K)$
  - $p(\boldsymbol{x}|\boldsymbol{z}=\boldsymbol{k}) = N(\boldsymbol{x}|\boldsymbol{\mu}_k,\boldsymbol{\Sigma}_k)$



• Clustering: The posterior p(z|x) identifies the mixture component

### **Example: Variational Autoencoder**



- A Mixture of an infinite number of Gaussians
  - $\boldsymbol{z} = N(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}), \, \boldsymbol{z} \in \mathbb{R}^h$
  - $p(\mathbf{x}|\mathbf{z}) = N(\mathbf{x}|\mu_{\theta}(\mathbf{z}), \Sigma_{\theta}(\mathbf{z}))$  where  $\mu_{\theta}, \Sigma_{\theta}$  are neural networks
  - Even though  $p(\mathbf{x}|\mathbf{z})$  is simple, the marginal  $p(\mathbf{x})$  is very complex/flexible
- Hope that after training, z will correspond to meaningful latent factors of variation (features)
- Unsupervised representation learning
- Features can be computed via  $p(\mathbf{z}|\mathbf{x})$

The Evidence Lower bound

$$\log p_{\theta}(\mathbf{x})$$

$$= E_{q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D\left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})\right)$$

$$+ D\left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}|\mathbf{x})\right)$$

og-likeliho

log  $p_{\theta}(x)$ ELBO  $KL(q_{\phi}(z), p_{\theta}(z|x))$ 

φ

- holds for  $\forall q_{\phi}(\boldsymbol{z}|\boldsymbol{x})$ , distribution of  $\boldsymbol{z}$  parameterized by  $\phi$  and dependent on  $\boldsymbol{x}$ . E.g.,  $q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) = N\left(\boldsymbol{z} \middle| \boldsymbol{\mu}_{\phi}(\boldsymbol{x}), \operatorname{diag}\left(\boldsymbol{\sigma}_{\phi}^{2}(\boldsymbol{x})\right)\right)$
- We want to jointly optimize over  $\theta$  and  $\phi$  to maximize the ELBO over a dataset D

### Variational Inference

• Variational inference requires that intractable posterior distributions be approximated by a class of known prob. dist.

 $E_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})] - D\left(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) \parallel p(\boldsymbol{z})\right)$ 

- We need to choose the computationally-feasible approximate posterior distribution  $q_{\phi}(\mathbf{z}|\mathbf{x})$
- We need efficient computation of

 $\nabla_{\boldsymbol{\phi}} E_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z})]$ 

### **Stochastic backpropagation**

Reparametrization

$$z = N(z|\mu, \operatorname{diag}(\sigma^2)) \Leftrightarrow z = \mu + \operatorname{diag}(\sigma^2)\epsilon, \quad \epsilon \sim N(0, I)$$

Backpropagation with Monte Carlo

 $\nabla_{\phi} E_{\boldsymbol{z} \sim q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}[f_{\theta}(\boldsymbol{z})] \Leftrightarrow E_{\boldsymbol{\epsilon} \sim N(\boldsymbol{\epsilon}|\boldsymbol{0},\boldsymbol{I})}[\nabla_{\phi} f_{\theta}(\boldsymbol{\mu}_{\phi} + \operatorname{diag}(\boldsymbol{\sigma}_{\phi}^{2})\boldsymbol{\epsilon})]$ 

### **Amortized Variational Inference**

- Inference network: model that learns an inverse map(encoder) from observations to latent variables
- Using this, we can compute a set of global variational parameters  $\phi$  valid for infrence at both training and test time
- The simplest inference models: diagonal Gaussian densities

$$q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) = N\left(\boldsymbol{z}|\boldsymbol{\mu}_{\phi}(\boldsymbol{x}), \operatorname{diag}\left(\boldsymbol{\sigma}_{\phi}^{2}(\boldsymbol{x})\right)\right)$$

### **VAE:** Autoencoder perspective

 $\mathcal{L}(\boldsymbol{x};\theta,\phi) = E_{q_{\phi}(\boldsymbol{z}|\boldsymbol{x})}[\log p_{\theta}(\boldsymbol{x}|\boldsymbol{z})] - D\left(q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) \parallel p(\boldsymbol{z})\right) \quad (\text{ELBO})$ 

- 1. Take a data point x', map it to sample  $\hat{z} \sim q_{\phi}(z|x')$  (encoder)
  - Sample  $\hat{z}$  from  $q_{\phi}(z|x') = N(z|\mu_{\phi}(x'), \text{diag}(\sigma_{\phi}^2(x')))$ , (encoder)
- 2. Reconstruct  $\hat{x}$  by sampling from  $p_{\theta}(x|\hat{z})$  (decoder or generator)
  - $\widehat{x} = G_{\theta}(\widehat{z})$
  - The training objective  $\mathcal{L}(\mathbf{x}; \theta, \phi)$ 
    - Reconstruction error
    - Regularizer

### The VAE training algorithm

- Sample minibatch of *n* training points  $x^{(1)}, x^{(2)}, \dots, x^{(n)}$  from  $p_{data}$
- Obtain *n* prob. distributions  $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$  using inference network
- Sample *n* latent feature points  $z^{(1)}, z^{(2)}, \dots, z^{(n)}$  from  $q_{\phi}(z|x^{(i)})$  resp.
- Update the generator parameters  $\phi$ ,  $\theta$  by stochastic gradient descent

$$\nabla_{\boldsymbol{\phi}} \ell(\boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{n} \nabla_{\boldsymbol{\phi}} \sum_{i=1}^{n} \left[ \left\| \boldsymbol{x}^{(i)} - G_{\boldsymbol{\theta}} \left( \boldsymbol{z}^{(i)} \right) \right\|_{2}^{2} - D \left( q_{\boldsymbol{\phi}} \left( \boldsymbol{z} | \boldsymbol{x}^{(i)} \right) \| p(\boldsymbol{z}) \right) \right]$$
$$\nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}, \boldsymbol{\phi}) = \frac{1}{n} \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{n} \left[ \left\| \boldsymbol{x}^{(i)} - G_{\boldsymbol{\theta}} \left( \boldsymbol{z}^{(i)} \right) \right\|_{2}^{2} \right]$$

• Repeat for fixed number of epochs

### Summary of Latent Variable Models

- Combine simple models to get a more flexible one (e.g., mixture of Gaussians)
- Directed model permits ancestral sampling (efficient generation):  $z \sim p(z), x \sim p_{\theta}(x|z)$
- However, log-likelihood is generally intractable. I.e., learning is difficult
- Joint learning of a model ( $\theta$ ) and an amortized inference component ( $\phi$ ) to achieve tractability via ELBO optimization
- Latent representations for any x can be inferred via  $q_{\phi}(z|x)$

### **Recap of likelihood-based learning**

- Model families:
  - Autoregressive models:  $p_{\theta}(\mathbf{x}) = \prod_{i=1}^{d} p_{\theta}(x_i | \mathbf{x}_{< i})$
  - Latent variable model:  $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z}$
- Autoregressive models provide tractable likelihoods but no direct mechanism for learning features
- Variational autoencoders can learn feature representations (via latent variables  $q_{\phi}(\mathbf{z}|\mathbf{x})$ ) but have intractable marginal likelihoods
- Key question: Can we design a latent variable model with tractable likelihoods?

### Normalizing flow models

- Consider a directed, latent variable model over observed variables *X* and latent variables *Z*
- In a normalizing flow model, the mapping between Z and X, given by  $f_{\theta} \colon \mathbb{R}^d \to \mathbb{R}^d$ , is deterministic and invertible such that  $X = f_{\theta}(Z)$  and  $Z = f_{\theta}^{-1}(X)$



### **Continuous random variables (recall)**

- If X is a continuous random variable, we can usually represent it using its probability density function  $p_X: \mathbb{R} \to \mathbb{R}^+$
- However, we cannot represent this function as a table anymore
- Typically consider parameterized densities:

• Gaussian: 
$$X = N(X|\mu, \sigma)$$
 if  $p_X(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-(x-\mu)^2/2\sigma^2}$ 

• Uniform: 
$$X = U(X|a, b)$$
 if  $p_X(x) = \frac{1}{b-a} \mathbb{1}_{[a \le x \le b]}$ 

• If *X* is a continuous random vector, we can usually represent it using its joint probability density function:

• Gaussian: 
$$p_X(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

### Change of variables formula

- Change of variables(1D case)
  - If X = f(Z) and  $f(\cdot)$  is monotone with inverse  $Z = f^{-1}(X)$ , then

$$p_X(x) = p_Z(f^{-1}(x))|(f^{-1})'(x)|$$

• Note that the "shape" of  $p_X(x)$  is different (more complex) from that of the prior  $p_Z(z)$ 

### Change of variables formula

- Change of variables(1D case):
  - If X = f(Z) and  $f(\cdot)$  is monotone with inverse  $Z = f^{-1}(X)$ , then

$$p_X(x) = p_Z(f^{-1}(x))|(f^{-1})'(x)|$$

Proof sketch: Assume f(·) is monotonically increasing
F<sub>X</sub>(x) = p[X ≤ x] = p[f(Z) ≤ x] = p[Z ≤ f<sup>-1</sup>(x)] = F<sub>Z</sub>(f<sup>-1</sup>(x))
Taking derivatives on both sides:

$$p_X(x) = \frac{dF_X(x)}{dx} = \frac{dF_Z(f^{-1}(x))}{dx} = p_Z(f^{-1}(x))(f^{-1})'(x)$$

- Recall from basic calculus that  $[f^{-1}]'(x) = \frac{1}{f'(f^{-1}(x))}$
- So, letting  $z = f^{-1}(x)$  we can also write

$$p_X(x) = p_Z(z) \left| \frac{1}{f'(z)} \right|$$

### **Geometry: Determinants and volumes**

- Let Z be a uniform random vector in  $[0,1]^d$
- Let X = AZ for a square invertible matrix A, with inverse  $A^{-1}$ . How is X distributed?
- Geometrically, the matrix A maps the unit hypercube  $[0,1]^d$  to a parallelotope
- Hypercube and parallelotope are generalizations of square/cube and parallelogram/parallelopiped to higher dimensions



• The matrix  $A = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$  maps a unit square to a parallelogram

### **Geometry: Determinants and volumes**

• The volume of the parallelotope is equal to the absolute value of the determinant of the matrix *A* 

$$\det(A) = \det\begin{pmatrix} a & c \\ b & d \end{pmatrix} = ad - bc$$



(a+c)(b+d) - ab - 2bc - cd = ad - bc

- Let X = AZ for a square invertible matrix A, with inverse  $A^{-1}$
- X is uniformly distributed over the parallelotope of area |det(A)|
- Hence, we have

 $p_X(\mathbf{x}) = p_Z(A^{-1}\mathbf{x})/|\det(A)| = p_Z(\mathbf{z})|\det(A^{-1})|$ 

Note similarity with 1D case formula

### Generalized change of variables

- For linear transformations specified via *A*, change in volume is given by the determinant of *A*
- For non-linear transformations  $f(\cdot)$ , the linearized change in volume is given by the determinant of the Jacobian of  $f(\cdot)$
- Change of variables (General case): The mapping between Z and X, given by  $f: \mathbb{R}^d \to \mathbb{R}^d$ , is invertible such that X = f(Z) and

$$Z = f^{-1}(X)$$

$$p_X(\mathbf{x}) = p_Z(\mathbf{f}^{-1}(\mathbf{x})) \left| \det\left(\frac{\partial \mathbf{f}^{-1}(\mathbf{x})}{\partial \mathbf{x}}\right) \right| = p_Z(\mathbf{z}) \left| \det\left(\frac{\partial \mathbf{f}(\mathbf{z})}{\partial \mathbf{z}}\right) \right|^{-1}$$

- Generalizes the 1D case  $p_X(x) = p_Z(f^{-1}(x))|(f^{-1})'(x)|$
- Unlike VAEs, x, z need to be continuous and have the same dimension

### **Two-Dimensional Example**

- Let  $Z = (Z_1, Z_2)^T$  be continuous random vector with joint density  $p_7$
- Let  $u: \mathbb{R}^2 \to \mathbb{R}^2$  be an invertible transformation denoted u = $(u_1, u_2)$  and let  $v = (v_1, v_2)$  be its inverse transformation
- Let  $X_1 = u_1(Z_1, Z_2)$  and  $X_2 = u_2(Z_1, Z_2)$
- Then,  $Z_1 = v_1(X_1, X_2)$  and  $Z_2 = v_2(X_1, X_2)$

 $\partial Z_1$ 

$$p_{X}(x_{1}, x_{2}) = p_{Z}(v_{1}(x_{1}, x_{2}), v_{2}(x_{1}, x_{2})) \left| \det \begin{pmatrix} \frac{\partial v_{1}(x_{1}, x_{2})}{\partial x_{1}} & \frac{\partial v_{1}(x_{1}, x_{2})}{\partial x_{2}} \\ \frac{\partial v_{2}(x_{1}, x_{2})}{\partial x_{1}} & \frac{\partial v_{2}(x_{1}, x_{2})}{\partial x_{2}} \end{pmatrix} \right|^{-1} \\ = p_{Z}(z_{1}, z_{2}) \left| \det \begin{pmatrix} \frac{\partial u_{1}(z_{1}, z_{2})}{\partial z_{1}} & \frac{\partial u_{1}(z_{1}, z_{2})}{\partial z_{2}} \\ \frac{\partial u_{2}(z_{1}, z_{2})}{\partial z_{2}} & \frac{\partial u_{2}(z_{1}, z_{2})}{\partial z_{2}} \end{pmatrix} \right|^{-1}$$

Deep Generative Models | mjgim@nims.re.kr |

 $\partial Z_2$ 

#### NIMS & AJOU University

### Normalizing flow models

- Consider a directed, latent variable model over observed variables *X* and latent variables *Z*
- In a normalizing flow model, the mapping between Z and X, given by  $f_{\theta} \colon \mathbb{R}^d \to \mathbb{R}^d$ , is deterministic and invertible such that  $X = f_{\theta}(Z)$  and  $Z = f_{\theta}^{-1}(X)$



### Normalizing flow models

• Using change of variables, the marginal likelihood is given by

$$p_X(\boldsymbol{x};\boldsymbol{\theta}) = p_Z\left(\boldsymbol{f}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{x})\right) \left| \det\left(\frac{\partial \boldsymbol{f}_{\boldsymbol{\theta}}^{-1}(\boldsymbol{x})}{\partial \boldsymbol{x}}\right) \right| = p_Z(\boldsymbol{z}) \left| \det\left(\frac{\partial \boldsymbol{f}_{\boldsymbol{\theta}}(\boldsymbol{z})}{\partial \boldsymbol{z}}\right) \right|^{-1}$$

• Note: *x*, *z* need to be continuous and have the same dimension



### **Normalizing Flows**

- A normalizing flow describes the transformation of a probability density through a sequence of invertible mappings
- By repeatedly applying the rule for change of variables, the initial density 'flows' through the sequence of invertible mappings
- Flexible, arbitrarily complex and scalable



### A Flow of Transformations

- **Normalizing**: Change of variables gives a normalized density after applying an invertible transformation
- Flow: Invertible transformations can be composed with each other

$$\boldsymbol{f}_{\theta}(\boldsymbol{z}_0) \coloneqq \boldsymbol{f}_K \circ \boldsymbol{f}_{K-1} \circ \cdots \circ \boldsymbol{f}_1(\boldsymbol{z}_0) = \boldsymbol{z}_K$$



### **A Flow of Transformations**



- Start with a simple distribution for  $z_0$  (e.g., Gaussian)
- Apply a sequence of *K* invertible transformations to finally obtain  $\mathbf{x} = \mathbf{z}_K$  $\mathbf{f}_{\theta}^{-1}(\mathbf{x}) = \mathbf{f}_1^{-1} \circ \mathbf{f}_2^{-1} \circ \cdots \circ \mathbf{f}_K^{-1}(\mathbf{x})$

### **A Flow of Transformations**

$$f_{\theta}(z_0) \coloneqq f_K \circ f_{K-1} \circ \cdots \circ f_1(z_0) = z_K = x$$
  
$$f_{\theta}^{-1}(x) = f_1^{-1} \circ f_2^{-1} \circ \cdots \circ f_K^{-1}(x)$$

• The marginal likelihood  $p_X(\mathbf{x})$  is given by

$$p_X(\mathbf{x};\theta) = p_Z\left(f_{\theta}^{-1}(\mathbf{x})\right) \left| \det\left(\frac{\partial f_{\theta}^{-1}(\mathbf{x})}{\partial \mathbf{x}}\right) \right|$$
$$= p_Z(\mathbf{z}) \left| \det\left(\frac{\partial f_{\theta}(\mathbf{z})}{\partial \mathbf{z}}\right) \right|^{-1}$$
$$= p_Z\left(f_{\theta}^{-1}(\mathbf{x})\right) \prod_{k=1}^{K} \left| \det\left(\frac{\partial f_k^{-1}(\mathbf{x}_k)}{\partial \mathbf{x}_k}\right) \right|$$

### Learning and Inference

• Learning via maximum likelihood over the dataset *D* 

$$\max_{\theta} \log p_X(D;\theta) = \sum_{\boldsymbol{x} \in D} \log p_Z\left(\boldsymbol{f}_{\theta}^{-1}(\boldsymbol{x})\right) + \log \left| \det\left(\frac{\partial \boldsymbol{f}_{\theta}^{-1}(\boldsymbol{x})}{\partial \boldsymbol{x}}\right) \right|$$

- Exact likelihood evaluation via inverse transformation  $x \mapsto z$  and change of variables formula
- Sampling via forward transformation  $z \mapsto x$

$$z \sim p_Z(z), x = f_\theta(z)$$

• Latent representations inferred via inverse transformation (no inference network required):  $z = f_{\theta}^{-1}(x)$ 

- Effect of normalizing flow on two distributions
- 10 planar transformations can transform simple distributions into a more complex one



• Planar flow

$$\boldsymbol{f}(\boldsymbol{z}) = \boldsymbol{z} + \boldsymbol{u}h(\boldsymbol{w}^T\boldsymbol{z} + b)$$

- parametrized by  $\theta = (w, u, b)$  where  $h(\cdot)$  is a nonlinear function
- Absolute value of the determinant of the Jacobian is given by

$$\left| \det\left(\frac{\partial \boldsymbol{f}(\boldsymbol{z})}{\partial \boldsymbol{z}}\right) \right| = \left| \det(\boldsymbol{I} + \boldsymbol{h}'(\boldsymbol{w}^T \boldsymbol{z} + \boldsymbol{b})\boldsymbol{u}\boldsymbol{w}^T) \right|$$
$$= \left| 1 + \boldsymbol{h}'(\boldsymbol{w}^T \boldsymbol{z} + \boldsymbol{b})\boldsymbol{u}^T \boldsymbol{w} \right| \quad (\text{matrix determinant lemma})$$

- Need to restrict parameters and non-linearity for the mapping to be invertible.
  - For example,  $h = \tanh(\cdot)$  and  $h'(\mathbf{w}^T \mathbf{z} + b)\mathbf{u}^T \mathbf{w} \ge -1$

• A family of transformations of Planar flows

$$\boldsymbol{f}_k(\boldsymbol{z}) = \boldsymbol{z} + \boldsymbol{u}_k h \big( \boldsymbol{w}_k^T \boldsymbol{z} + \boldsymbol{b}_k \big)$$

- parametrized by  $(w_k, u_k, b_k)$  where  $h(\cdot)$  is a nonlinear function
- Let *f*<sub>θ</sub>(*z*<sub>0</sub>) ≔ *f*<sub>K</sub> ∘ … *f*<sub>2</sub> ∘ *f*<sub>1</sub>(*z*<sub>0</sub>) = *z*<sub>K</sub> = *x*The log-likelihood *p*<sub>X</sub>(*x*; θ) is given by

$$\log p_X(\boldsymbol{x};\theta) = \log p_Z(\boldsymbol{z}) - \sum_{k=1}^K \log |1 + h' (\boldsymbol{w}_k^T \boldsymbol{z}_{k-1} + b) \boldsymbol{u}_k^T \boldsymbol{w}_k|$$





Approximating four non-Gaussian 2D distributions



(d) Comparison of KL-divergences.

### Desiderata for flow models

- Simple prior  $p_Z(z)$  that allows for efficient sampling and tractable likelihood evaluation. E.g., isotropic Gaussian
- Invertible transformations with tractable evaluation:
  - Likelihood evaluation requires efficient evaluation of  $x \mapsto z$ mapping
  - Sampling requires efficient evaluation of  $z \mapsto x$  mapping

### **Desiderata for flow models**

- Computing likelihoods also requires the evaluation of determinants  $d \times d$  of Jacobian matrices
- Computing the determinant for a  $d \times d$  matrix is  $O(d^3)$ 
  - Key idea: Choose transformations so that the resulting Jacobian matrix has special structure
  - For example, the determinant of a triangular matrix is the product of the diagonal entries, i.e., an O(d) operation

### **Triangular Jacobian**

$$\boldsymbol{x} = (x_1, x_2, \cdots, x_d)^T = \boldsymbol{f}(\boldsymbol{z}) = \left(f_1(\boldsymbol{z}), f_2(\boldsymbol{z}), \cdots, f_d(\boldsymbol{z})\right)^T$$
$$\frac{\partial \boldsymbol{f}}{\partial \boldsymbol{z}} = \begin{pmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & \frac{\partial f_1}{\partial z_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_d}{\partial z_1} & \cdots & \frac{\partial f_d}{\partial z_d} \end{pmatrix}$$
If  $x_i = f_i(\boldsymbol{z})$  only depends on  $\boldsymbol{z}_{\leq i}$ , then
$$\frac{\partial \boldsymbol{f}}{\partial \boldsymbol{z}} = \begin{pmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial f_d}{\partial z_1} & \cdots & \frac{\partial f_d}{\partial z_d} \end{pmatrix}$$

• has lower triangular structure

### Recap of normalizing flow models

- Transform simple to complex distributions via sequence of invertible transformations
- Directed latent variable models with marginal likelihood given by the change of variables formula
- Triangular Jacobian permits efficient evaluation of log-likelihood
- Examples
  - Invertible transformations with diagonal Jacobians (NICE, Real-NVP)
  - Autoregressive Models as Normalizing Flow Models
  - Invertible CNNs (MintNet)
  - Gaussianization flows

### **Designing invertible transformations**

- Nonlinear Independent Components Estimation (Dinh et al., 2014) composes two kinds of invertible transformations: additive coupling layers and rescaling layers
- Real-NVP (Dinh et al., 2017)
- Inverse Autoregressive Flow (Kingma et al., 2016)
- Masked Autoregressive Flow (Papamakarios et al., 2017)
- I-resnet (Behrmann et al, 2018)
- Glow (Kingma et al, 2018)
- MintNet (Song et al., 2019)
- • •

### **NICE - Additive coupling layers**

- Partition the variables z into two disjoint subsets, say  $z_{\leq h}$  and  $z_{>h}$  for any  $1 \leq h < d$
- Forward mapping  $z \mapsto x$ :
  - $x_{\leq h} = z_{\leq h}$  (identity transformation)
  - $x_{>h} = z_{>h} + m_{\theta}(z_{\leq h})$   $(m_{\theta}(\cdot)$  is a neural network with parameters  $\theta$ , h input units, and d h output units)
- Inverse mapping  $x \mapsto z$ :
  - $\mathbf{z}_{\leq h} = \mathbf{x}_{\leq h}$  (identity transformation)
  - $\mathbf{z}_{>h} = \mathbf{x}_{>h} m_{\theta}(\mathbf{x}_{\leq h})$

### **NICE - Additive coupling layers**

• Jacobian of forward mapping:

$$\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}} = \begin{pmatrix} I_h & \boldsymbol{0} \\ \frac{\partial \boldsymbol{x}_{>h}}{\partial \boldsymbol{z}_{\le h}} & I_{d-h} \end{pmatrix}$$
$$\det\left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right) = 1$$

• Volume preserving transformation since determinant is 1

### NICE - Rescaling layers

- Additive coupling layers are composed together (with arbitrary partitions of variables in each layer)
- Final layer of NICE applies a rescaling transformation
- Forward mapping  $z \mapsto x$ :

 $x_i = s_i z_i$ 

- where  $s_i > 0$  is the scaling factor for the i<sup>th</sup> dimension
- Inverse mapping  $x \mapsto z$ :

$$z_i = \frac{x_i}{s_i}$$

• Jacobian of forward mapping

 $\operatorname{diag}(s_1, s_2, \cdots, s_d)$ 

### Samples generated via NICE



(a) Model trained on MNIST

(b) Model trained on TFD

### Samples generated via NICE



(c) Model trained on SVHN

(d) Model trained on CIFAR-10

### **Real-NVP: Non-volume preserving extension**

- Forward mapping  $z \mapsto x$ :
  - $x_{\leq h} = z_{\leq h}$  (identity transformation)
  - $\mathbf{x}_{>h} = \mathbf{z}_{>h} \odot \exp(\alpha_{\theta}(\mathbf{z}_{\leq h})) + \mu_{\theta}(\mathbf{z}_{\leq h})$
  - $\alpha_{\theta}(\cdot)$  and  $\mu_{\theta}(\cdot)$  are both neural networks with parameters  $\theta$ , *h* input units, and d - h output units
  - $\odot$  denotes elementwise product
- Inverse mapping  $x \mapsto z$ :
  - $\mathbf{z}_{\leq h} = \mathbf{x}_{\leq h}$  (identity transformation)
  - $\mathbf{z}_{>h} = (\mathbf{x}_{>h} \mu_{\theta}(\mathbf{x}_{\le h})) \odot \exp(-\alpha_{\theta}(\mathbf{x}_{\le h}))$

### **Real-NVP: Non-volume preserving extension**

• Jacobian of forward mapping:

$$\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}} = \begin{pmatrix} l_h & \boldsymbol{0} \\ \frac{\partial \boldsymbol{x}_{>h}}{\partial \boldsymbol{z}_{\le h}} & \text{diag}(\exp(\alpha_\theta(\boldsymbol{z}_{\le h}))) \end{pmatrix}$$
$$\det\left(\frac{\partial \boldsymbol{x}}{\partial \boldsymbol{z}}\right) = \prod_{i=h+1}^d \exp(\alpha_\theta(\boldsymbol{z}_{\le h})_i) = \exp\left(\sum_{i=h+1}^d \alpha_\theta(\boldsymbol{z}_{\le h})_i\right)$$

• Non-volume preserving transformation in general since determinant can be less than or greater than 1

### Samples generated via Real-NVP



### **Continuous Autoregressive models as flow models**

• Consider a Gaussian autoregressive model:

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^{d} p(x_i | \mathbf{x}_{< i})$$

• such that  $p(x_i | \mathbf{x}_{< i}) = N(x_i | \mu_i(\mathbf{x}_{< i}), \exp(2\alpha_i(\mathbf{x}_{< i})))$ . Here  $\mu_i$  and  $\alpha_i$  are neural networks for i > 1 and  $\mu_1$  and  $\alpha_1$  are constants

### **Continuous Autoregressive models as flow models**

• Consider a Gaussian autoregressive model:

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^{d} p(x_i | \mathbf{x}_{< i})$$

- such that  $p(x_i | \mathbf{x}_{< i}) = N(x_i | \mu_i(\mathbf{x}_{< i}), \exp(2\alpha_i(\mathbf{x}_{< i})))$ . Here  $\mu_i$  and  $\alpha_i$  are neural networks for i > 1 and  $\mu_1$  and  $\alpha_1$  are constants
- Sampler for this model:
  - Sample  $z_i \sim N(0,1)$  for  $i = 1, \dots, d$
  - Let  $x_1 = \exp(\alpha_1) z_1 + \mu_1$ . Compute  $\mu_2(x_1), \alpha_2(x_1)$
  - Let  $x_2 = \exp(\alpha_2) z_2 + \mu_2$ . Compute  $\mu_3(x_1, x_2), \alpha_3(x_1, x_2)$
  - Let  $x_3 = \exp(\alpha_3) z_3 + \mu_3$ . ...
- Flow interpretation: transforms samples from the standard Gaussian z to those generated from the model x via invertible transformations (parameterized by  $\mu_i$  and  $\alpha_i$ )

### Masked Autoregressive Flow (MAF)



- Forward mapping  $z \mapsto x$ :
  - Let  $x_1 = \exp(\alpha_1) z_1 + \mu_1$ . Compute  $\mu_2(x_1), \alpha_2(x_1)$
  - Let  $x_2 = \exp(\alpha_2) z_2 + \mu_2$ . Compute  $\mu_3(x_1, x_2), \alpha_3(x_1, x_2)$
  - Let  $x_3 = \exp(\alpha_3) z_3 + \mu_3$ . ...
- Sampling is sequential and slow (like autoregressive): O(d) time

### Masked Autoregressive Flow (MAF)



- **Inverse** mapping  $x \mapsto z$ :
  - Compute all  $\mu_i$ ,  $\alpha_i$
  - Let  $z_1 = (x_1 \mu_1) / \exp(\alpha_1)$  (scale and shift)
  - Let  $z_2 = (x_2 \mu_2) / \exp(\alpha_2)$
  - Let  $z_3 = (x_3 \mu_3) / \exp(\alpha_3) \dots$
- Jacobian is lower diagonal, hence efficient determinant computation
- Likelihood evaluation is easy and parallelizable

### **Inverse Autoregressive Flow (IAF)**



- Forward mapping  $z \mapsto x$  (parallel):
  - Sample  $z_i \sim N(0,1)$  for  $i = 1, \dots, d$
  - Compute all  $\mu_i$ ,  $\alpha_i$
  - Let  $x_1 = \exp(\alpha_1) z_1 + \mu_1$
  - Let  $x_2 = \exp(\alpha_2) z_2 + \mu_2$
  - Let  $x_3 = \exp(\alpha_3) z_3 + \mu_3 \cdots$

### **Inverse Autoregressive Flow (IAF)**



- **Inverse** mapping  $x \mapsto z$  (sequential):
  - Let  $z_1 = (x_1 \mu_1) / \exp(\alpha_1)$ . Compute  $\mu_2(z_1), \alpha_2(z_1)$
  - Let  $z_2 = (x_2 \mu_2) / \exp(\alpha_2)$ . Compute  $\mu_3(z_1, z_2), \alpha_3(z_1, z_2)$
  - Let  $z_3 = (x_3 \mu_3) / \exp(\alpha_3)$ . ...
- Fast to sample
- Slow to evaluate likelihoods of data points
- Note: Fast to evaluate likelihoods of a generated point

### **MAF and IAF**



- IAF is an inverse of MAF
- Interchanging z and x in the inverse transformation of MAF gives the forward transformation of IAF
- Similarly, forward transformation of MAF is inverse transformation of IAF

### **MAF and IAF**



- Computational tradeoffs
  - MAF: Fast likelihood evaluation, slow sampling
  - IAF: Fast sampling, slow likelihood evaluation
- MAF more suited for training based on MLE, density estimation
- IAF more suited for real-time generation

### Summary of Normalizing flow models

- Transform simple distributions into more complex distributions via change of variables
- Jacobian of transformations should have tractable determinant for efficient learning and density estimation
- Computational tradeoffs in evaluating forward and inverse transformations



## Thanks