# Normalizing Flow Models

#### Stefano Ermon, Aditya Grover

Stanford University

Lecture 7

Stefano Ermon, Aditya Grover (AI Lab)

Deep Generative Models

Lecture 7 1 / 21

# Recap of likelihood-based learning so far:



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

# Simple Prior to Complex Data Distributions

- Desirable properties of any model distribution:
  - Analytic density
  - Easy-to-sample
- Many simple distributions satisfy the above properties e.g., Gaussian, uniform distributions
- Unfortunately, data distributions could be much more complex (multi-modal)
- Key idea: Map simple distributions (easy to sample and evaluate densities) to complex distributions (learned via data) using change of variables.

- Let Z be a uniform random variable  $\mathcal{U}[0,2]$  with density  $p_Z$ . What is  $p_Z(1)$ ?  $\frac{1}{2}$
- Let X = 4Z, and let  $p_X$  be its density. What is  $p_X(4)$ ?
- $p_X(4) = p(X = 4) = p(4Z = 4) = p(Z = 1) = p_Z(1) = 1/2$  No
- Clearly, X is uniform in [0, 8], so  $p_X(4) = 1/8$

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

$$p_X(x) = p_Z(h(x))|h'(x)|$$

- Previous example: If X = 4Z and  $Z \sim \mathcal{U}[0, 2]$ , what is  $p_X(4)$ ?
- Note that h(X) = X/4
- $p_X(4) = p_Z(1)h'(4) = 1/2 \times 1/4 = 1/8$

### Geometry: Determinants and volumes

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



• The volume of the parallelotope is equal to the determinant of the transformation *A* 

$$\det(A) = \det \left(\begin{array}{cc} a & c \\ b & d \end{array}\right) = ad - bc$$

• X is uniformly distributed over the parallelotope. Hence, we have

$$p_X(\mathbf{x}) = p_Z(W\mathbf{x}) |\det(W)|$$
$$= p_Z(W\mathbf{x}) / |\det(A)|$$

## 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(·), the *linearized* change in volume is given by the determinant of the Jacobian of f(·).
- Change of variables (General case): The mapping between Z and X, given by  $\mathbf{f} : \mathbb{R}^n \mapsto \mathbb{R}^n$ , is invertible such that  $X = \mathbf{f}(Z)$  and  $Z = \mathbf{f}^{-1}(X)$ .

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

- Note 1: **x**, **z** need to be continuous and have the same dimension. For example, if  $\mathbf{x} \in \mathbb{R}^n$  then  $\mathbf{z} \in \mathbb{R}^n$
- Note 2: For any invertible matrix A,  $det(A^{-1}) = det(A)^{-1}$

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

#### Two Dimensional Example

- Let  $Z_1$  and  $Z_2$  be continuous random variables with joint density  $p_{Z_1,Z_2}$ .
- Let  $u = (u_1, u_2)$  be a transformation
- Let  $v = (v_1, v_2)$  be the 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  $Z_2 = v_2(X_1, X_2)$

$$p_{X_1,X_2}(x_1,x_2)$$

$$= p_{Z_1,Z_2}(v_1(x_1,x_2),v_2(x_1,x_2)) \left| \det \left( \begin{array}{c} \frac{\partial v_1(x_1,x_2)}{\partial x_1} & \frac{\partial v_1(x_1,x_2)}{\partial x_2} \\ \frac{\partial v_2(x_1,x_2)}{\partial x_1} & \frac{\partial v_2(x_1,x_2)}{\partial x_2} \end{array} \right) \right| \text{(inverse)}$$

$$= p_{Z_1,Z_2}(z_1,z_2) \left| \det \left( \begin{array}{c} \frac{\partial u_1(z_1,z_2)}{\partial z_1} & \frac{\partial u_1(z_1,z_2)}{\partial z_2} \\ \frac{\partial u_2(z_1,z_2)}{\partial z_1} & \frac{\partial u_2(z_1,z_2)}{\partial z_2} \end{array} \right) \right|^{-1} \text{(forward)}$$

# 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  $\mathbf{f}_{\theta} : \mathbb{R}^{n} \mapsto \mathbb{R}^{n}$ , is deterministic and invertible such that  $X = \mathbf{f}_{\theta}(Z)$ and  $Z = \mathbf{f}_{\theta}^{-1}(X)$



• Using change of variables, the marginal likelihood  $p(\mathbf{x})$  is given by

$$p_X(\mathbf{x}; heta) = p_Z\left(\mathbf{f}_{ heta}^{-1}(\mathbf{x})
ight) \left| \det\left(rac{\partial \mathbf{f}_{ heta}^{-1}(\mathbf{x})}{\partial \mathbf{x}}
ight) 
ight|$$

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

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

$$\mathbf{x} \triangleq \mathbf{z}_{M} = \mathbf{f}_{\theta}^{M} \circ \cdots \circ \mathbf{f}_{\theta}^{1}(\mathbf{z}_{0}) = \mathbf{f}_{\theta}^{M}(\mathbf{f}_{\theta}^{M-1}(\cdots(\mathbf{f}_{\theta}^{1}(\mathbf{z}_{0})))) \triangleq \mathbf{f}_{\theta}(\mathbf{z}_{0})$$

- Start with a simple distribution for  $z_0$  (e.g., Gaussian)
- Apply a sequence of *M* invertible transformations

$$p_X(\mathbf{x}; heta) = p_Z\left(\mathbf{f}_{ heta}^{-1}(\mathbf{x})
ight) \prod_{m=1}^M \left|\det\left(rac{\partial(\mathbf{f}_{ heta}^m)^{-1}}{\partial \mathbf{z}_m}
ight)
ight|$$

(determininant of product equals product of determinants)

### Planar flows

• Planar flow (Rezende & Mohamed, 2016). Invertible transformation

$$\mathbf{x} = \mathbf{f}_{\theta}(\mathbf{z}) = \mathbf{z} + \mathbf{u}h(\mathbf{w}^{\mathsf{T}}\mathbf{z} + b)$$

parameterized by  $\theta = (\mathbf{w}, \mathbf{u}, b)$  where  $h(\cdot)$  is a non-linearity

• Absolute value of the determinant of the Jacobian is given by

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

 Need to restrict parameters and non-linearity for the mapping to be invertible. For example, h = tanh() and h'(w<sup>T</sup>z + b)u<sup>T</sup>w ≥ -1

## Planar flows

• Base distribution: Gaussian



• Base distribution: Uniform



• 10 planar transformations can transform simple distributions into a more complex one

Stefano Ermon, Aditya Grover (AI Lab)

### Learning and Inference

• Learning via maximum likelihood over the dataset  ${\cal D}$ 

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

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

$$\mathbf{z} \sim p_Z(\mathbf{z}) \ \mathbf{x} = \mathbf{f}_{ heta}(\mathbf{z})$$

• Latent representations inferred via inverse transformation (no inference network required!)

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

- Simple prior  $p_Z(\mathbf{z})$  that allows for efficient sampling and tractable likelihood evaluation. E.g., isotropic Gaussian
- Invertible transformations with tractable evaluation:
  - $\bullet\,$  Likelihood evaluation requires efficient evaluation of  $x\mapsto z$  mapping
  - $\bullet$  Sampling requires efficient evaluation of  $z\mapsto x$  mapping
- Computing likelihoods also requires the evaluation of determinants of  $n \times n$  Jacobian matrices, where *n* is the data dimensionality
  - Computing the determinant for an  $n \times n$  matrix is  $O(n^3)$ : prohibitively expensive within a learning loop!
  - Key idea: Choose tranformations 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(n) operation

### Triangular Jacobian

$$\mathbf{x} = (x_1, \cdots, x_n) = \mathbf{f}(\mathbf{z}) = (f_1(\mathbf{z}), \cdots, f_n(\mathbf{z}))$$

$$J = \frac{\partial \mathbf{f}}{\partial \mathbf{z}} = \begin{pmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & \frac{\partial f_1}{\partial z_n} \\ \cdots & \cdots & \cdots \\ \frac{\partial f_n}{\partial z_1} & \cdots & \frac{\partial f_n}{\partial z_n} \end{pmatrix}$$

Suppose  $x_i = f_i(\mathbf{z})$  only depends on  $\mathbf{z}_{\leq i}$ . Then

$$J = \frac{\partial \mathbf{f}}{\partial \mathbf{z}} = \begin{pmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & 0\\ \cdots & \cdots & \cdots\\ \frac{\partial f_n}{\partial z_1} & \cdots & \frac{\partial f_n}{\partial z_n} \end{pmatrix}$$

has lower triangular structure. Determinant can be computed in **linear time**. Similarly, the Jacobian is upper triangular if  $x_i$  only depends on  $\mathbf{z}_{>i}$ 

- NICE or 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)

# NICE - Additive coupling layers

Partition the variables  ${\bf z}$  into two disjoint subsets, say  ${\bf z}_{1:d}$  and  ${\bf z}_{d+1:n}$  for any  $1 \leq d < n$ 

- Forward mapping  $\mathbf{z} \mapsto \mathbf{x}$ :
  - $\mathbf{x}_{1:d} = \mathbf{z}_{1:d}$  (identity transformation)
  - $\mathbf{x}_{d+1:n} = \mathbf{z}_{d+1:n} + m_{\theta}(\mathbf{z}_{1:d}) (m_{\theta}(\cdot) \text{ is a neural network with parameters}$ 
    - heta, d input units, and n-d output units)
- Inverse mapping  $\mathbf{x} \mapsto \mathbf{z}$ :
  - $\mathbf{z}_{1:d} = \mathbf{x}_{1:d}$  (identity transformation)
  - $z_{d+1:n} = x_{d+1:n} m_{\theta}(x_{1:d})$
- Jacobian of forward mapping:

$$J = \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \begin{pmatrix} I_d & 0\\ \frac{\partial \mathbf{x}_{d+1:n}}{\partial \mathbf{z}_{1:d}} & I_{n-d} \end{pmatrix}$$
$$\det(J) = 1$$

#### • 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 → x:

$$x_i = s_i z_i$$

where  $s_i > 0$  is the scaling factor for the *i*-th dimension.

• Inverse mapping  $\mathbf{x} \mapsto \mathbf{z}$ :

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

• Jacobian of forward mapping:

$$J = \mathsf{diag}(\mathbf{s})$$

$$\det(J) = \prod_{i=1}^n s_i$$



(a) Model trained on MNIST

(b) Model trained on TFD



(c) Model trained on SVHN

(d) Model trained on CIFAR-10