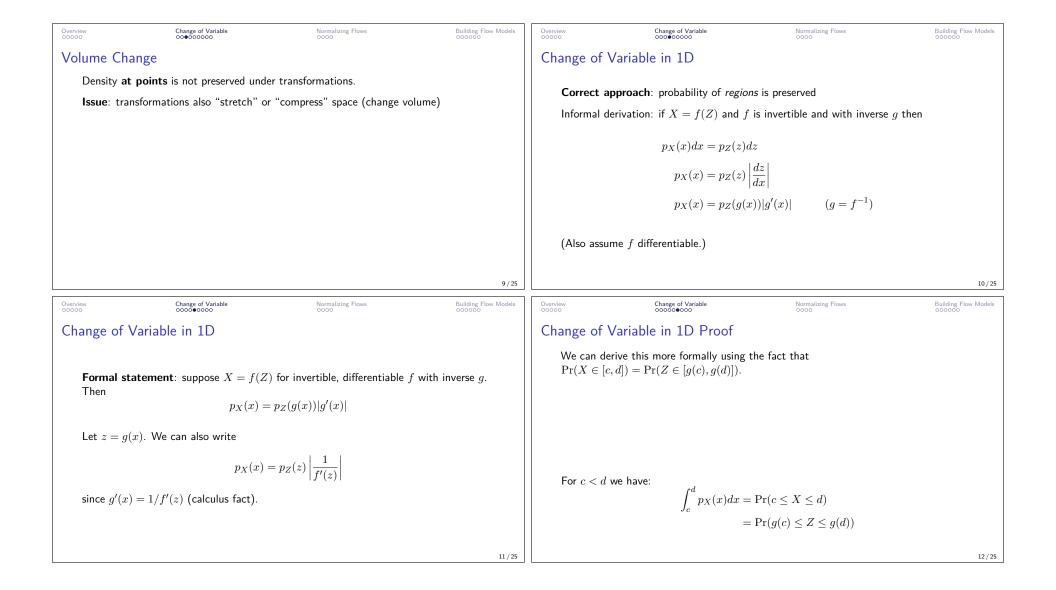
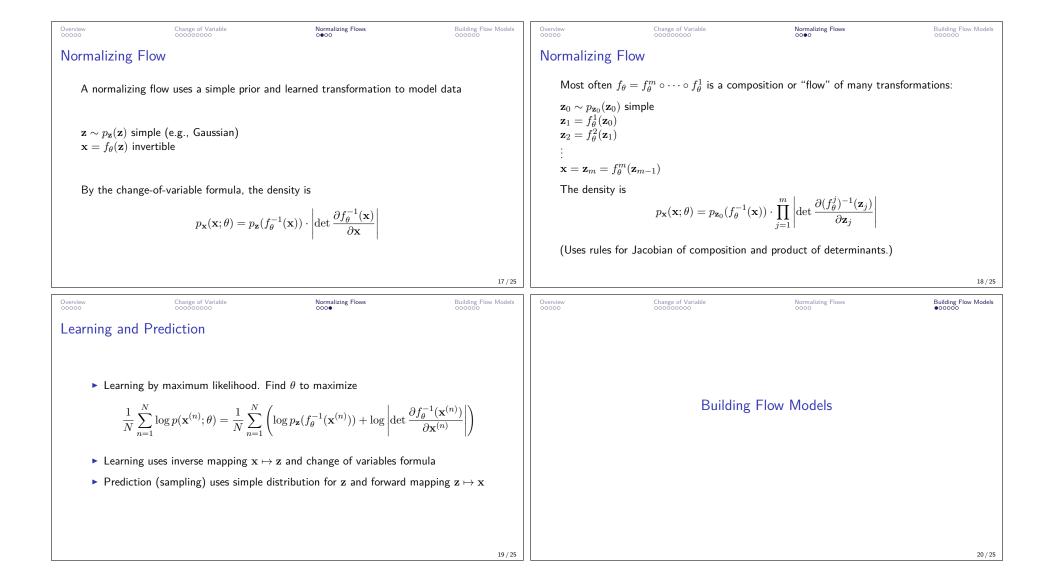
COMPSCI 688: Probabilistic Graphical Models Lecture 23: Normalizing Flows Overview Dan Sheldon Maning College of Information and Computer Sciences University of Massachusetts Andrest Overview Method Massachusetts Markett 1/30 Coverview Network Interference of Information and Computer Sciences University of Massachusetts Andrest 1/30 Coverview Network Inform (a) Simple Distribution 1/30 Cover Overview Motivation: Transforming a Simple Distribution Suppose we want to learn a model $p_0(x)$ for a complex x (like images). What properties do we want from $p_0(x)$? Easy to sample (useful for generation) Consider our VAE model $p_0(x)$ but with no noise $x \sim p(x)$ simple, e.g. $\mathcal{N}(0, 1)$ $x = f_0(x)$ $y = p_0(x)$ But data distributions are complex [E.g. multi-modal. Key idea bashind flow models: may simple distributions to complex ones through deterministic invertible transformations Could we learn $p_0(x)$ directify by MLE? Could we learn $p_0(x)$ directify by MLE? • Can easily generate samples $x \sim p_0(x)$ To learn, med to compute the density $p_0(x)$ under transformation f_0 . Can we do it? Demo	Overview 00000	Change of Variable 00000000	Normalizing Flows	Building Flow Models	Overview •0000	Change of Variable 00000000	Normalizing Flows	Building Flow Models		
Dan Sheldon Manning College of Information and Computer Sciences University of Massachusetts Amherst Petially based on materials by Benjamit M. Martin (mutitelly unmass eth) and Justin Danke (denoteds unmass eth) 1/25 Overciew Charge of Variable Motivation: Transforming a Simple Distribution Suppose we want to learn a model $p_0(x)$ for a complex x (like images). What properties do we want from $p_0(x)$? • Easy to sample (useful for generation) • Easy to sample distributions satisfy these properties (e.g., Gaussian, uniform). But data distributions are complex! E.g. multi-modal. Key idea behind flow models: may simple distributions to complex ones through			1	lels						
Image: distribution is statisfy these properties (e.g., Gaussian, uniform). I/25 Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (ike images). What properties do we want from $p_{\theta}(x)$? Image: distribution is a complex t (e.g., Gaussian, uniform). But data distributions are complex! E.g. multi-modal. Key idea behind flow models: may simple distributions to complex ones through the tensitivity in tensin tensity $p_{\theta}(x)$ under transformation f_{θ} . Can we do it?		Da	n Sheldon			(Overview			
$\begin{array}{ c c c c c c } \hline 1/25 \end{array}$										
Motivation: Transforming a Simple DistributionMotivation: Transforming a Simple Distribution (Learning)Suppose we want to learn a model $p_{\theta}(\mathbf{x})$ for a complex \mathbf{x} (like images). What properties do we want from $p_{\theta}(\mathbf{x})$?Consider our VAE model $p_{\theta}(\mathbf{x})$ but with no noise $\mathbf{z} \sim p(\mathbf{z})$ simple, e.g. $\mathcal{N}(0, I)$ $\mathbf{x} = f_{\theta}(\mathbf{z})$ • Easy to sample (useful for generation) • Easy to evaluate density (useful for learning)Consider our VAE model $p_{\theta}(\mathbf{x})$ but with no noise $\mathbf{x} = f_{\theta}(\mathbf{z})$ Many simple distributions satisfy these properties (e.g., Gaussian, uniform). But data distributions are complex! E.g. multi-modal. Key idea behind flow models: map simple distributions to complex ones through deterministic invertibilis invertibilis to repreference.Could we learn $p_{\theta}(\mathbf{x})$ "directly" by MLE? • Can easily generate samples $\mathbf{x} \sim p_{\theta}(\mathbf{x})$ • To learn, need to compute the density $p_{\theta}(\mathbf{x})$ under transformation f_{θ} . Can we do it?		Partially based on materials by Benjamin M. Marlin	marlin@cs.umass.edu) and Justin Domke (domke@c	,				2/25		
Suppose we want to learn a model $p_{\theta}(\mathbf{x})$ for a complex \mathbf{x} (like images). What properties do we want from $p_{\theta}(\mathbf{x})$? • Easy to sample (useful for generation) • Easy to evaluate density (useful for learning) Many <i>simple</i> distributions satisfy these properties (e.g., Gaussian, uniform). But data distributions are <i>complex</i> ! E.g. multi-modal. Key idea behind flow models: map simple distributions to complex ones through deterministic invertible temeformations	Overview oeooo	Change of Variable	Normalizing Flows	Building Flow Models	Overview 00000	Change of Variable	Normalizing Flows	Building Flow Models		
do we want from $p_{\theta}(\mathbf{x})$? • Easy to sample (useful for generation) • Easy to evaluate density (useful for learning) Many simple distributions satisfy these properties (e.g., Gaussian, uniform). But data distributions are complex! E.g. multi-modal. Key idea behind flow models: map simple distributions to complex ones through deterministic investible transformation f_{θ} . Can we do it?	Motivatio	n: Transforming a Simpl	e Distribution		Motivation: Transforming a Simple Distribution (Learning)					
But data distributions are <i>complex</i> ! E.g. multi-modal. Key idea behind flow models : map simple distributions to complex ones through deterministic investible twentime $p_{\theta}(\mathbf{x})$ under transformation f_{θ} . Can we do it?	 do we want from p_θ(x)? ► Easy to sample (useful for generation) 				$\mathbf{z} \sim p(\mathbf{z})$ simple, e.g. $\mathcal{N}(0, I)$ \implies $p_{\theta}(\mathbf{x})$					
Key idea behind flow models: map simple distributions to complex ones through f_{θ} . Can we do it? To learn, need to compute the density $p_{\theta}(\mathbf{x})$ under transformation f_{θ} . Can we do it?	Many <i>simple</i> distributions satisfy these properties (e.g., Gaussian, uniform).					Could we learn $p_{\theta}(\mathbf{x})$ "directly" by MLE?				
3/25	Key idea behind flow models: map simple distributions to complex ones through				► To lea			$f_{ heta}.$ Can we do it?		

00000	Change of Variable	Normalizing Flows	Building Flow Models	Overview 0000	Change of Variable	Normalizing Flows	Building Flow Models	
Motivation:	Transforming a Simple	e Distribution (Inferenc	e)	Can we do i	t?			
Also useful	in variational inference, e.g.	$q_{\phi}(\mathbf{z})$ in VAEs		Not in gen	eral. Consider the VAE mode	I		
Goal: $q_{\phi}(\mathbf{z}$	$\mathbf{z}) pprox p(\mathbf{z} \mathbf{x})$ where p , \mathbf{x} are given by	ven. We used reparameterized	d Gaussians:					
					• ()	$\mathcal{N}(\mathbf{z};0,I)$ ("easy")		
	$\frac{\epsilon \sim \mathcal{N}(0, I)}{\mathbf{z} = \mathcal{T}_{\ell}(\epsilon) = L\epsilon + u} =$	$\Rightarrow \mathbf{z} \sim q_{\phi}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mu, L)$	L^T)		$\mathbf{x} \sim p(\mathbf{x} \mathbf{z})$			
What if we	e used complex $\mathcal{T}_{\phi}(\epsilon)$ (e.g. net				gh $p(\mathbf{z})$ is "easy", $p(\mathbf{x}) = \int p(\mathbf{x}) d\mathbf{x}$ have produced \mathbf{x} .	$(\mathbf{z})p(\mathbf{x} \mathbf{z})d\mathbf{z}$ is "hard": need to	o enumerate all ${f z}$	
► Would	have a rich class of variation	al distributions.		Even if x =	= $f_{ heta}(\mathbf{z})$ is deterministic, could	be hard to reason about ${f z}$ the	hat produced \mathbf{x} .	
	easily sample from $q_{\phi}(\mathbf{z})$	· · · (-) · · · · · · · · · · · · · · · · · · ·	T Can un da it?	But if f_{θ} is invertible , we can do it!				
FOR EL	LBO, need to compute density	$q_{\phi}(\mathbf{z})$ under transformation	T_{ϕ} . Can we do it?					
			5 / 25				6 /	
hiomious	Change of Variable	Normalizing Flows	Building Flow Models	Overview	Change of Variable	Normalizing Flows	Building Flow Mode	
)0000	€00Õ0000	0000	000000	Overview 00000	0000000	0000	000000	
0000	€0000000	0000	000000		/ariable in 1D (False St		000000	
VVEL VIEW	€0000000	0000	000000	Change of $\$			000000	
Vectore Vectore	●oocoooo	0000	000000	Change of \backslash Example ($Z \sim Unif(0)$	/ariable in 1D (False St false start). Suppose 0,1)		000000	
			000000	Change of V Example ($Z \sim \text{Unif}(C)$ X = 2Z + C	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
7464 MEW 20000		e of Variable	000000	Change of \backslash Example ($Z \sim Unif(0)$	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
900000			00000	Change of V Example ($Z \sim \text{Unif}(C)$ X = 2Z + C	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
20000 20000			00000	Change of V Example ($Z \sim \text{Unif}(C)$ X = 2Z + C	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
90000			00000	Change of V Example ($Z \sim \text{Unif}(C)$ X = 2Z + C	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
OOOOO			00000	Change of V Example ($Z \sim \text{Unif}(C)$ X = 2Z + C	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)		000000	
Overview 00000			00000	Change of V Example ($Z \sim \text{Unif}(0)$ X = 2Z + What is p_X	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z) f(2)?	cart)	00000	
OOOOO			00000	Change of V Example ($Z \sim \text{Unif}(0)$ X = 2Z + What is p_X Easy to gu	/ariable in 1D (False St false start). Suppose (0,1) 1 := f(Z)	cart) $_{Z}(\frac{1}{2}) = 1.$ Wrong.	000000	



Overview 00000	Change of Variable	Normalizing Flows	Building Flow Models	Overview 00000	Change of Variable 0000000€0	Normalizing Flows	Building Flow Models		
	$=\int_{a}^{g(d)}$	$p_Z(z)dz$		Change of	Variable: General Case				
$= \int_{g(c)}^{g(d)} p_Z(z) dz$ $= \int_c^d p_Z(g(x))g'(x) dx$ The last line uses the calculus change of variable formula with the substitution $z = g(x)$. (So this is really the same change of variable formula.) Since c and d are arbitrary, by comparing the integrands we see that $p_X(x) = p_Z(g(x))g'(x)$.					Suppose $\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})$ and $x = f(\mathbf{z})$ for invertible, differentiable $f : \mathbb{R}^D \to \mathbb{R}^D$ with inverse g . Then $p_x(\mathbf{x}) = p_{\mathbf{z}}(g(\mathbf{x})) \cdot \left \det \frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \right $ • The matrix $\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} \in \mathbb{R}^{D \times D}$ is the Jacobian of g . It's (i, j) th entry is $\frac{\partial g_i(x)}{\partial x_j}$ • It's also true that $\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}} = \left(\frac{\partial f(\mathbf{z})}{\partial \mathbf{z}}\right)^{-1}$ for $\mathbf{z} = g(\mathbf{x})$. So we often call $\frac{\partial g(\mathbf{x})}{\partial \mathbf{x}}$ the inverse Jacobian of f				
			13/25				14 / 25		
Overview 00000	Change of Variable 0000000●	Normalizing Flows	Building Flow Models	Overview 00000	Change of Variable	Normalizing Flows ●000	Building Flow Models		
• Another version, often convenient. Let $\mathbf{z} = g(\mathbf{x})$. Then $p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) \cdot \left \det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} \right ^{-1}$ • Geometrically, $\left \det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} \right $ describes how much f changes the volume of a small hypercube.					Norma	lizing Flows			
			15 / 25				16 / 25		



Overview 00000	Change of Variable	Normalizing Flows	Building Flow Models	Overview 00000	Change of Variable 00000000	Normalizing Flows	Building Flow Models			
Building Flow Models					Triangular Jacobian					
	a flow model we need tribution $p(\mathbf{z})$ that is "easy". C	an sample and compute den	sity.	$J = \frac{\partial f}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & \frac{\partial f_1}{\partial z_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D}{\partial z_1} & \cdots & \frac{\partial f_D}{\partial z_D} \end{bmatrix}$						
Trans	sformations $f_ heta$ that are			Suppose <i>x</i>	$x_i = f_i(\mathbf{z})$ only depends on z_1	$,\ldots,z_i$. Then				
 Always invertible Allow us to compute the determinant easily. In general, it is O(D³) — too expensive! Key idea: choose tranformations with special structure 					$J = \frac{\partial f}{\partial \mathbf{z}} = \begin{bmatrix} \frac{\partial f_1}{\partial z_1} & \cdots & 0\\ \vdots & \ddots & \vdots\\ \frac{\partial f_D}{\partial z_1} & \cdots & \frac{\partial f_D}{\partial z_D} \end{bmatrix}$					
					iangular \implies the determinan ted in linear time .	t is the product of the diagon:	al entries of J , can			
			21 / 25				22 / 25			
Overview 00000	Change of Variable	Normalizing Flows 0000	Building Flow Models	Overview 00000	Change of Variable	Normalizing Flows	Building Flow Models 0000€0			
Real-NVP				The inverse mapping is $\mathbf{x}\mapsto \mathbf{z}$ is therefore						
	many constructions that ensur ". We split z and x into two e	0		$ \begin{aligned} \mathbf{z}_1 &= \mathbf{x}_1 & \text{(identity)} \\ \mathbf{z}_2 &= (\mathbf{x}_2 - \mu_\theta(\mathbf{x}_1)) \odot \exp(-\alpha_\theta(\mathbf{x}_1)) & \text{(unshift and unscale } \mathbf{x}_2 \text{ based on } \mathbf{x}_1 \text{)} \end{aligned} $						
$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}.$ The forward mapping $\mathbf{z} \mapsto \mathbf{x}$ is					The Jacobian of the forward mapping and its determinant are $J = \frac{\partial \mathbf{x}}{\partial \mathbf{z}} = \begin{bmatrix} I & 0\\ \frac{\partial \mathbf{x}_2}{\partial \mathbf{z}_1} & diag(\exp(\alpha_\theta(\mathbf{z}_1)) \end{bmatrix}$					
										$ \begin{aligned} \mathbf{x}_1 &= \mathbf{z}_1 & (\text{identity}) \\ \mathbf{x}_2 &= \mu_{\theta}(\mathbf{z}_1) + \mathbf{z}_2 \odot \exp(\alpha_{\theta}(\mathbf{z}_1)) & (\text{shift and scale } \mathbf{z}_2 \text{ based on } \mathbf{z}_1) \end{aligned} $
where $\mu_{ heta}($	(\cdot) and $lpha_ heta(\cdot)$ are neural network	is from $\mathbb{R}^d \to \mathbb{R}^d$.		Change or	der of dimensions in different	layers, so sometimes $\mathbf{z}_2\mapsto\mathbf{x}_2$	is identity instead.			
				1						

Overview 00000	Change of Variable	Normalizing Flows	Building Flow Models				
Demo							
•	Demo: implementation and 2d densi	ty estimation with Real-NVP					
•	There are tons of examples on the ir look.	nternet of images generated by flow	vs. Take a				
•	Flows have been used for tons of ap	plications					
	 They can be extremely good for VI. They are good at generating images, but not the most competitive models right now (if you care). One reason is they restrict f_θ too much. Some more competitive current models descend from normalizing flows. 						
			25 / 25				