### Invertible Normalizing Flows

ECE57000: Artificial Intelligence

David I. Inouye

#### GAN Limitation:

Cannot compute density values

- Evaluation of GANs is challenging
  - Explicit density models could use test log likelihood
  - "I think this looks better than that"
  - Inception-based scores

- Cannot use for classification or outlier detection
- Normalizing flows provide exact density values



# Common problem with GANs: Mode collapse hinders diversity of samples

From: https://developers.google.com/machine-learning/gan/problems

### Normalizing flows do not suffer from mode collapse as MLE is used



Metz, L., Poole, B., Pfau, D., & Sohl-Dickstein, J. (2016). Unrolled generative adversarial networks. *arXiv preprint arXiv:1611.02163*.

(f) True Data (g) GAN http://papers.nips.cc/paper/6923-veegan-reducing-modecollapse-in-gans-using-implicit-variational-learning.pdf



https://software.intel.com/en-us/blogs/2017/08/21/mode-collapse-in-gans

GAN Limitation: Cannot go from observed to latent space, i.e.  $x \rightarrow z$  not possible/easy

- Cannot manipulate an observed image in latent space
  - Cannot do the following,  $x \to z$ , z' = z + 3,  $z' \to x'$
  - Rather, must start from fake image based on random



Z



# Normalizing flows enable interpolation between real images



Figure 5: Linear interpolation in latent space between real images.

https://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions.pdf

### Normalizing flows enable transformations of **real image** along various features



(a) Smiling

(b) Pale Skin



(c) Blond Hair





(e) Young

(f) Male

Figure 6: Manipulation of attributes of a face. Each row is made by interpolating the latent code of an image along a vector corresponding to the attribute, with the middle image being the original image

#### https://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions.pdf

Normalizing flows enable more powerful inference distributions  $q_f(z|x)$  in VAEs

- The probabilistic encoder in VAEs requires 2 things:
  - 1. Ability to sample from  $q_f(z|x)$  via reparameterization trick
  - 2. Ability to compute exact density of  $q_f(z|x)$  for  $KL(q_f(z|x), p_g(z)) = \mathbb{E}_{z \sim q_f(z|x)} \left[ \log \frac{q_f(z|x)}{p_g(z)} \right]$
- Normalizing flows have these capabilities and thus significantly generalize the Gaussian q<sub>f</sub>(z|x)

Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling, M. (2016). Improved variational inference with inverse autoregressive flow. *Advances in neural information processing systems*, *29*, 4743-4751.

z space

x space

 $p_q(x|z)$ 

 $q_f(z|x)$ 

#### Overview of Normalizing Flows



7

Normalizing flows use invertible deep models for the generator which allow more capabilities

- Transforming between observed/input and latent space is easy
  - x = G(z)
  - $\blacktriangleright z = G^{-1}(x)$
- Simple sampling like GANs
  - $z \sim$  SimpleDistribution
  - $x = G(z) \sim \hat{p}_g(x)$ , which is estimated distribution
- Exact density is computable via change of variables
   Standard maximum likelihood estimation can be used for training

### Comparing VAEs and normalizing flows: Flows give zero reconstruction error



### Comparing GANs and normalizing flows: Normalizing flows can use MLE training



Output

Back to maximum likelihood estimation (MLE): <u>How</u> can we compute the likelihood for normalizing flows?

- Suppose
  - *z* ~ Uniform([0,1]), i. e., *p<sub>z</sub>(z)* = 1 (latent space is uniform)
  - G(z) = 2z

• Thus, 
$$x = G(z) = 2z$$
.

What is the density function of x, what is p<sub>x</sub>(x)? <u>Change of variables formula</u> gives  $p_x$  in terms of the  $p_z$  and the derivative of  $G^{-1}$ 

- Key idea: Must conserve density volume (so that distribution sums to 1).
- $p_x(x)|dx| = p_z(z)|dz|$ , this is like the preservation of volume/area/mass.
  - Intuition: We only have 1 unit of "dirt" to move around.

• Rearrange above equation to get formula  $p_x(x) = \left| \frac{dz}{dx} \right| p_z(z) = \left| \frac{dG^{-1}(x)}{dx} \right| p_z(G^{-1}(x))$ 

#### Demo of change of variables

### Derivation of change of variables using CDF function (Increasing)

Assume x = G(z), where G(z) is an increasing function, i.e.,  $z_1 \le z_2 \Rightarrow G(z_1) \le G(z_2)$  $z_1 \le z_2 \Rightarrow G^{-1}(z_1) \le G^{-1}(z_2)$ 

Remember CDF

$$F_x(a) = \Pr(x \le a) = \int_{-\infty}^{\infty} p_x(t) dt$$

ra.

Derivation of change of variables using CDF function (Increasing)

- From the previous slide, we have that  $F_x(a) = F_z(G^{-1}(a))$
- Now take the derivative of both sides with respect to a
- $\frac{dF_{x}(a)}{da} = \frac{dF_{z}(G^{-1}(a))}{da}$   $p_{x}(a) = \frac{dF_{z}(G^{-1}(a))}{d(G^{-1}(a))} \left(\frac{dG^{-1}(a)}{da}\right)$   $p_{x}(a) = p_{z}(G^{-1}(a)) \left(\frac{dG^{-1}(a)}{da}\right)$   $p_{x}(a) = p_{z}(G^{-1}(a)) \left(\frac{dG^{-1}(a)}{da}\right)$
- What about decreasing functions?

### Derivation of change of variables using CDF function (Decreasing)

Assume x = G(z), where G(z) is a decreasing function, i.e.,  

$$z_1 \leq z_2 \Rightarrow G(z_1) \geq G(z_2)$$
  
 $z_1 \leq z_2 \Rightarrow G^{-1}(z_1) \geq G^{-1}(z_2)$ 
F<sub>x</sub>(a) = Pr(x ≤ a) = Pr(G(z) ≤ a)
= Pr(G^{-1}(G(z)) \leq G^{-1}(a))
= Pr(z \geq G^{-1}(a))
= 1 - F\_z(G^{-1}(a))

Now take the derivative of both sides with respect to a

$$\frac{dF_{x}(a)}{da} = p_{x}(a)$$
$$-\frac{dF_{z}(G^{-1}(a))}{da} = -\frac{dF_{z}(G^{-1}(a))}{d(G^{-1}(a))} \left(\frac{dG^{-1}(a)}{da}\right)$$
$$= -p_{z}(G^{-1}(a)) \left(\frac{dG^{-1}(a)}{da}\right) = p_{z}(G^{-1}(a)) \left|\frac{dG^{-1}(a)}{da}\right|$$

- ►  $z \sim \text{Uniform}([0,1])$
- $v \sim$  AnotherDistribution
- $x = F_v^{-1}(z)$ , where  $F_v^{-1}$  is the inverse CDF for v
- What is the distribution of x?

• 
$$p_x(x) = p_z(F_v(x)) \left| \frac{dF_v(x)}{dx} \right|$$
  
•  $p_x(x) = (1)|p_v(x)| = p_v(x)$ 



### What about change of variables in higher dimensions?

- Let's again build a little intuition (see demo)
- Again, conservation of volume: Consider infinitesimal expansion or shrinkage of volume p(x<sub>1</sub>, x<sub>2</sub>)|dx<sub>1</sub>dx<sub>2</sub>| = p(z<sub>1</sub>, z<sub>2</sub>)|dz<sub>1</sub>dz<sub>2</sub>|
- Given that Jacobian is all mixed derivatives we get generalization for vector to vector invertible functions:

$$p_x(x) = |\det J_{G^{-1}}(x)| p_z(G^{-1}(x))$$

Interpretation: What is the Jacobian again? The best linear approximation at a point



The determinant measures the local *linear* expansion or shrinkage around a point



Fact: The determinant Jacobian of compositions of functions is the product of determinant Jacobians

- Suppose  $F(x) = F_2(F_1(x))$
- The Jacobian expands like the chain rule  $J_F(x) = J_{F_2}(F_1(x))J_{F_1}(x) = J_{F_2}J_{F_1}$
- If we take the determinant of the Jacobian, then it becomes a product of determinants

$$\det J_F = \det J_{F_2} J_{F_1} = (\det J_{F_2}) (\det J_{F_1})$$

This will be useful since each layer of our flows will be invertible Okay, now back to learning flows: The log likelihood is the sum of determinant terms for each layer

Simply optimize the minimize negative log likelihood where  $F_{\theta} = G^{-1}$  $\arg\min_{F_{\theta}} -\frac{1}{n} \sum_{i} \log p_{x}(x_{i};\theta)$  $-\frac{1}{n} \sum_{i} \log p_{z} \left( F_{\theta}(x_{i}) \right) \left| \det J_{F_{\theta}}(x_{i}) \right|$  $-\frac{1}{n} \sum_{i} \left[ \log p_z (F_\theta(x_i)) + \log \left| \det J_{F_\theta}(x_i) \right| \right]$  $-\frac{1}{n} \sum_{i} \left[ \log p_z \left( F_\theta(x_i) \right) + \sum_{\ell} \log \left| \det J_{F_0^{(\ell)}} \left( z_i^{(\ell-1)} \right) \right| \right]$ where  $z_i^0 = x_i$ , and  $z_i^{\ell} = F_{\Delta}^{(\ell)}(z_i^{\ell-1})$ 

### Overview of Normalizing Flows

22



Normalizing flow architectures: How do we create these invertible layers?

- Consider arbitrary invertible transformation  $F_{\theta}$ 
  - How often would  $\left|\det J_{F_{\theta}}\right|$  need to be computed?
- High computation costs
  - Determinant costs roughly O(d<sup>3</sup>) even if Jacobian is already computed!
  - Would need to be computed every stochastic gradient iteration

How do we create these invertible layers? Independent transformation on each dimension

► 
$$z_1 = F_1(x_1)$$
  
►  $z_2 = F_2(x_2)$   
►  $z_3 = F_3(x_3)$ 

What is the Jacobian?

$$J_F = \begin{bmatrix} \frac{dF_1(x_1)}{dx_1} & 0 & 0 \\ 0 & \frac{dF_2(x_2)}{dx_2} & 0 \\ 0 & 0 & \frac{dF_3(x_3)}{dx_3} \end{bmatrix}$$

### How do we create these invertible layers? <u>Autoregressive Flows</u> based on chain rule

- Forward Density estimation (in parallel)
  - $\blacktriangleright z_1 = F_1(x_1)$
  - $rightarrow z_2 = F_2(x_2|x_1)$
  - $\bullet z_3 = F_3(x_3 | x_1, x_2)$
- Inverse Sampling (conditioned on x so must be sequential)
  - $x_1 = F_1^{-1}(z_1)$
  - $x_2 = F_2^{-1}(z_2 | x_1)$
  - $\bullet x_3 = F_3^{-1}(z_3 | x_1, x_2)$
- What is the Jacobian and determinant?
  - Product of diagonal!

Rezende, D., & Mohamed, S. (2015, June). Variational Inference with Normalizing Flows. In *International Conference on Machine Learning* (pp. 1530-1538).

$$J_F = \begin{bmatrix} \frac{dF_1}{dx_1} & 0 & 0\\ \frac{dF_2}{dx_1} & \frac{dF_2}{dx_2} & 0\\ \frac{dF_3}{dx_1} & \frac{dF_3}{dx_2} & \frac{dF_3}{dx_3} \end{bmatrix}$$

### How do we create these invertible layers? Inverse Autoregressive Flows based on chain rule

- Forward Density estimation (sequential)
  - $\blacktriangleright z_1 = F_1(x_1)$
  - $rightarrow z_2 = F_2(x_2 | z_1)$
  - $rightarrow z_3 = F_3(x_3 | z_1, z_2)$
- Inverse Sampling (parallel)
  - $x_1 = F_1^{-1}(z_1)$  $x_2 = F_2^{-1}(z_2|z_1)$
  - $x_3 = F_3^{-1}(z_3|z_1, z_2)$
- What is the Jacobian and determinant?
  - Product of diagonal!

Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., & Welling, M. (2016). Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems* (pp. 4743-4751).

$$J_F = \begin{bmatrix} \frac{dF_1}{dx_1} & 0 & 0\\ \frac{dF_2}{dx_1} & \frac{dF_2}{dx_2} & 0\\ \frac{dF_3}{dx_1} & \frac{dF_3}{dx_2} & \frac{dF_3}{dx_3} \end{bmatrix}$$

Scale-and-shift simple form of invertible functions (MAF <u>https://arxiv.org/pdf/1705.07057.pdf</u>)

Forward – Density estimation (parallel)

$$z_1 = \exp(\alpha_1)x_1 + \mu_1$$
  

$$z_2 = \exp(\alpha_2)x_2 + \mu_2, \ \alpha_2 = f_2(x_1), \ \mu_2 = g_2(x_1)$$
  

$$z_3 = \exp(\alpha_3)x_3 + \mu_3, \ \alpha_3 = f_3(x_1, x_2), \ \mu_3 = g_3(x_1, x_2)$$

What is the Jacobian and determinant?

$$J_F = \begin{bmatrix} \exp(\alpha_1) & 0 & 0 \\ \frac{dz_2}{dx_1} & \exp(\alpha_2) & 0 \\ \frac{dz_3}{dx_1} & \frac{dz_3}{dx_2} & \exp(\alpha_3) \end{bmatrix}$$

## RealNVP and GLOW: Several key architecture ideas for **image-based** normalizing flows

- 1. Coupling layers
- 2. Invertible squeeze operation
- 3. Split dimensions along channel
- 4. Hierarchical structure
- 5. 1 x 1 convolutions

<u>Coupling layers</u> allow parallel density estimation **and** sampling

Keep some set of features fixed and transform others

$$Z_{1:i-1} = x_{1:i-1}$$
  
$$Z_{i:d} = \exp(f(x_{1:i-1})) \odot x_{i:d} + g(x_{1:i-1})$$

- Reverse or shuffle coordinates and repeat
- What is Jacobian?

$$J_F = \begin{bmatrix} I & 0\\ J_{cross} & \text{diag}(\exp(f(x_{1:i-1}))) \end{bmatrix}$$

Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density estimation using real nvp. arXiv preprint arXiv:1605.08803.

How to split dimensions for coupling layers? The squeeze operation trades off between spatial and channel dimensions



### $H/2 \times W/2 \times 4C$

Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density estimation using real nvp. arXiv preprint arXiv:1605.08803.

HxWxC

How to split dimensions for coupling layers? Checkboard or channel-wise masking can be used to separate fixed and non-fixed set of variables

White are **fixed**, i.e.,  $x_{1:i-1}$ , and black are **transformed**,  $x_{i:d}$ .





Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density estimation using real nvp. arXiv preprint arXiv:1605.08803.

Hierarchical factorization is like an invertible dimensionality reduction method

- After each block, half of the dimensions are fixed and the rest pass through more transformations
- Intuitively, the important part of the signal propagates deeper



Dinh, L., Sohl-Dickstein, J., & Bengio, S. (2016). Density estimation using real nvp. arXiv preprint arXiv:1605.08803.

GLOW: Convolutional flows 1 x 1 invertible convolutions are like fully connected layers for each pixel

- Image tensor:  $h \times w \times c$
- If we use c filters than we map from a  $\mathbf{h}\times w\times c$  to another  $\mathbf{h}\times w\times c$  image
- The number of parameters is a matrix  $W \in \mathbb{R}^{c \times c}$
- 1x1 convolutions can be seen as a linear transform along the channel dimension (mixes the channel dimensions)



# Highly realistic random samples from powerful flow model (GLOW)



Figure 1: Synthetic celebrities sampled from our model; see Section 3 for architecture and method, and Section 5 for more results.

https://papers.nips.cc/paper/8224-glow-generative-flow-with-invertible-1x1-convolutions.pdf

### More recent normalizing flows

- Neural Spline Flows and Neural Autoregressive Flows – Add more flexible 1D transforms
- Flow++ Careful tweaks to some previous models
- FFJORD Uses neural Ordinary Differential Equations (ODE) to implicitly define invertible functions
- Residual Flows Careful construction of residual networks that are invertible (uses Lipschitz idea)
- MaCow Masked Convolutional Generative Flow (carefully constructed masked convolutions to ensure invertibility)

Similar concepts can be used to generate realistic audio (WaveGlow)

Listen to some examples <u>https://nv-adlr.github.io/WaveGlow</u>

Very similar concepts for audio generation