

DSC291: Machine Learning with Few Labels

Deep generative modeling
Generative adversarial learning

Zhiting Hu

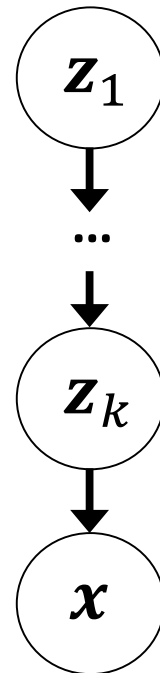
Lecture 14, February 10, 2023

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

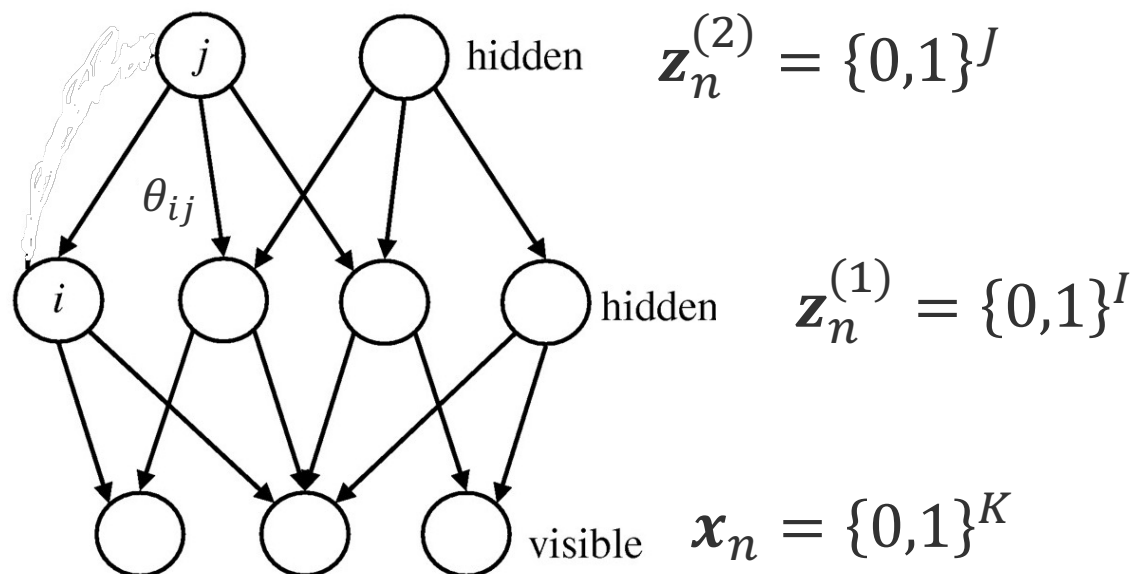
Deep generative models

- Define probabilistic distributions over a set of variables
- "Deep" means multiple layers of hidden variables!



Early forms of deep generative models

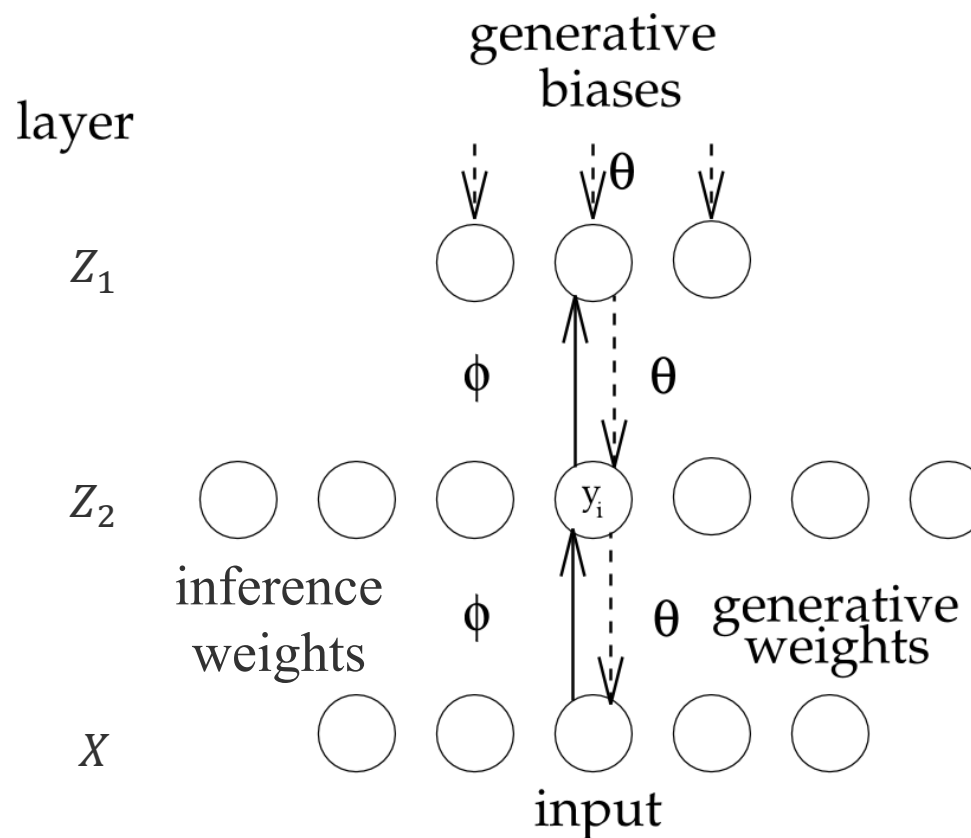
- Hierarchical Bayesian models
 - Sigmoid belief nets [Neal 1992]



$$p\left(x_{kn} = 1 \mid \boldsymbol{\theta}_k, \mathbf{z}_n^{(1)}\right) = \sigma\left(\boldsymbol{\theta}_k^T \mathbf{z}_n^{(1)}\right)$$
$$p\left(z_{in}^{(1)} = 1 \mid \boldsymbol{\theta}_i, \mathbf{z}_n^{(2)}\right) = \sigma\left(\boldsymbol{\theta}_i^T \mathbf{z}_n^{(2)}\right)$$

Early forms of deep generative models

- Hierarchical Bayesian models
 - Sigmoid belief nets [Neal 1992]
- Neural network models
 - Helmholtz machines [Dayan et al.,1995]



[Dayan et al. 1995]

Early forms of deep generative models

- Hierarchical Bayesian models
 - Sigmoid belief nets [Neal 1992]
- Neural network models
 - Helmholtz machines [Dayan et al.,1995]
 - Predictability minimization [Schmidhuber 1995]

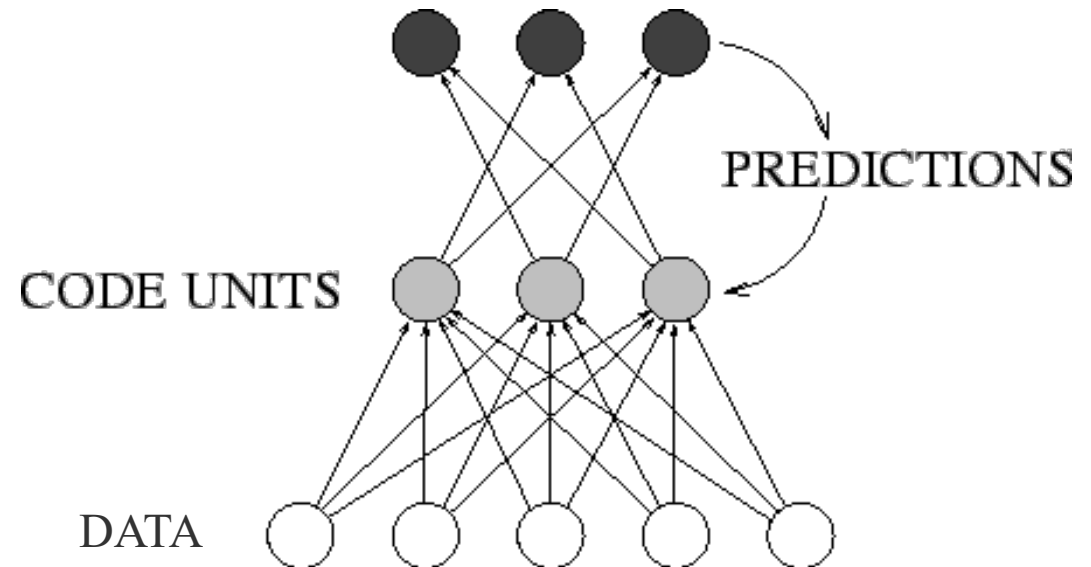


Figure courtesy: Schmidhuber 1996

Early forms of deep generative models

- Training of DGMs via an EM style framework

- Sampling / data augmentation

$$\mathbf{z} = \{\mathbf{z}_1, \mathbf{z}_2\}$$

$$\mathbf{z}_1^{new} \sim p(\mathbf{z}_1 | \mathbf{z}_2, \mathbf{x})$$

$$\mathbf{z}_2^{new} \sim p(\mathbf{z}_2 | \mathbf{z}_1^{new}, \mathbf{x})$$

- Variational inference

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}, \mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x}) || p(\mathbf{z})) := \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$

$$\max_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{x})$$

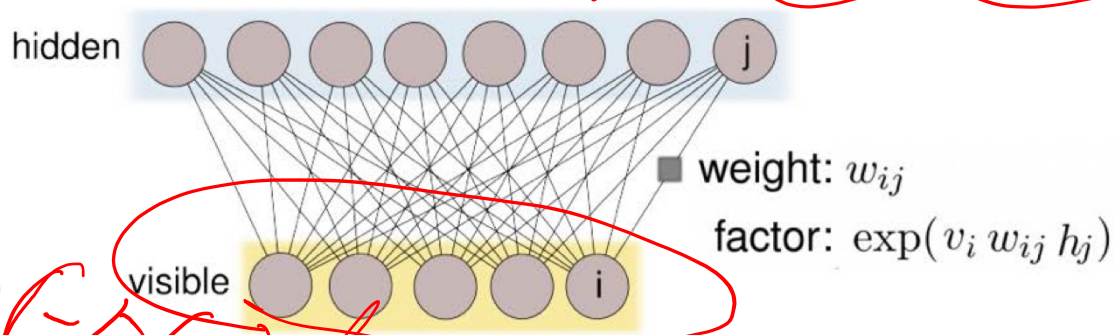
- Wake sleep

$$\text{Wake: } \min_{\boldsymbol{\theta}} \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})]$$

$$\text{Sleep: } \min_{\boldsymbol{\phi}} \mathbb{E}_{p_\theta(\mathbf{x}|\mathbf{z})}[\log q_\phi(\mathbf{z}|\mathbf{x})]$$

Resurgence of deep generative models

- Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
 - Building blocks of deep probabilistic models



$-E(v, h)$

$\exp\left\{\sum_{ij} v_i w_{ij} h_j\right\}$

$p(v, h) =$

energy-based

$Z = \int \exp(-E(x)) dx$

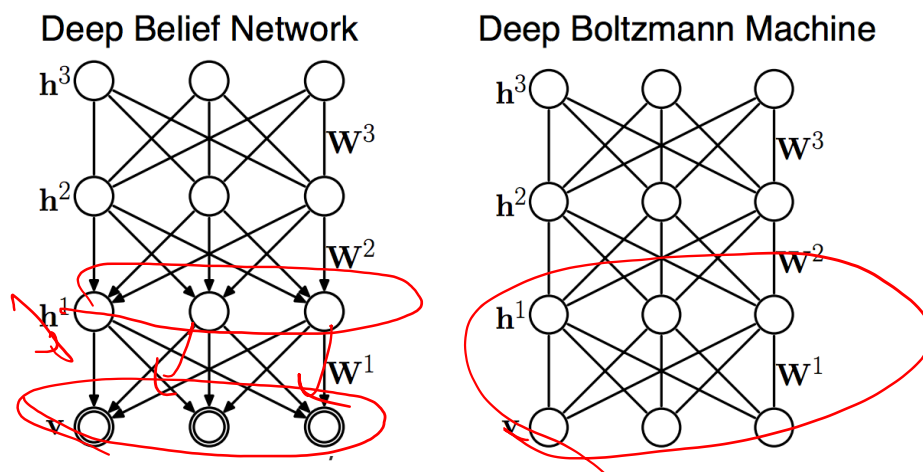
$p(x) = \frac{\exp(-E(x))}{Z}$

fun ans

Caribbs distributions

Resurgence of deep generative models

- Restricted Boltzmann machines (RBMs) [Smolensky, 1986]
 - Building blocks of deep probabilistic models
- Deep belief networks (DBNs) [Hinton et al., 2006]
 - Hybrid graphical model
 - Inference in DBNs is problematic due to explaining away
- Deep Boltzmann Machines (DBMs) [Salakhutdinov & Hinton, 2009]
 - Undirected model



Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]

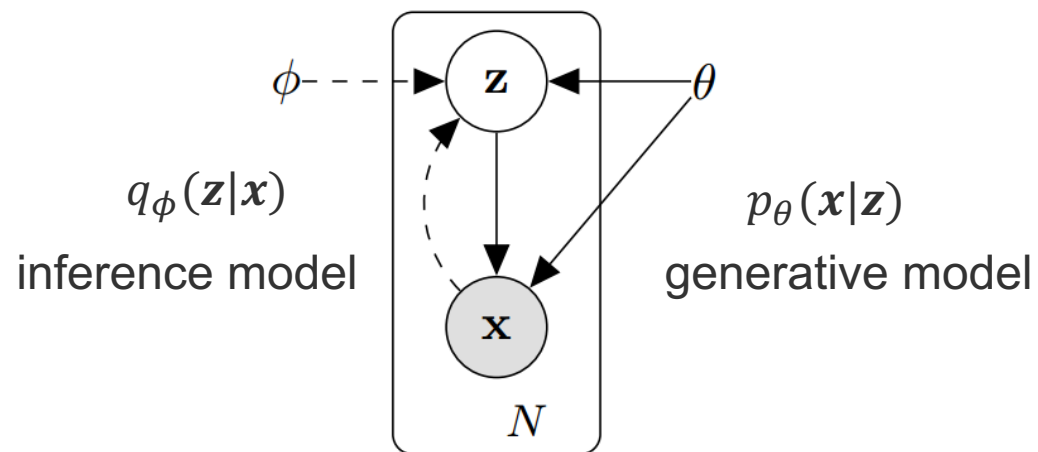
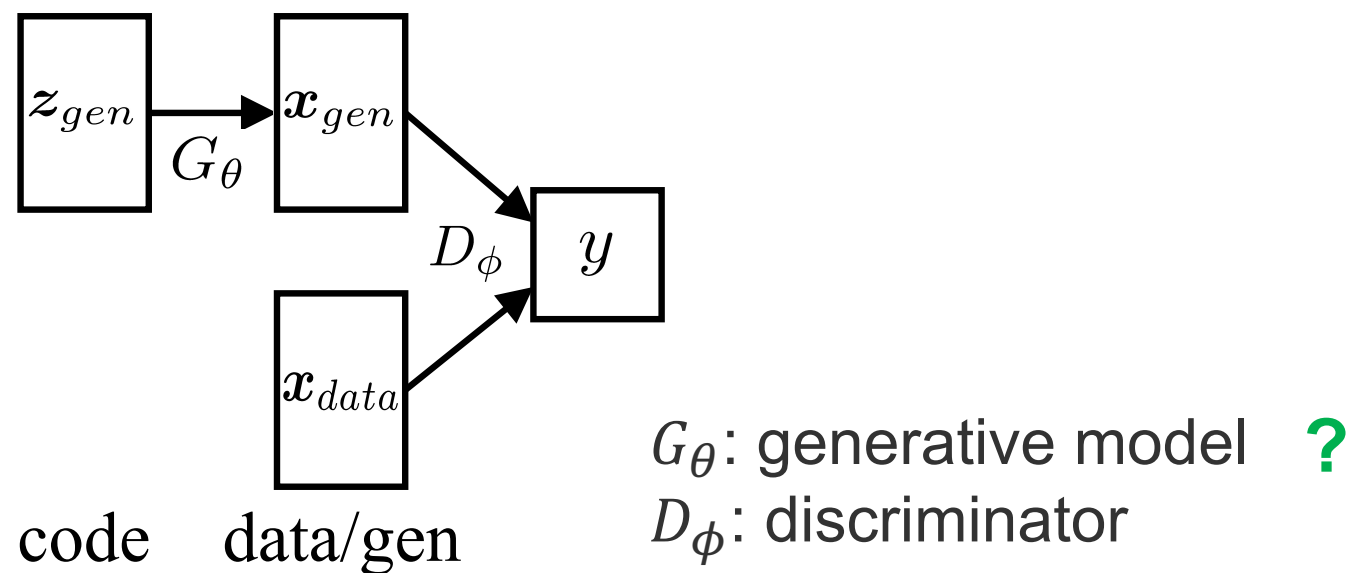


Figure courtesy: Kingma & Welling, 2014

Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]

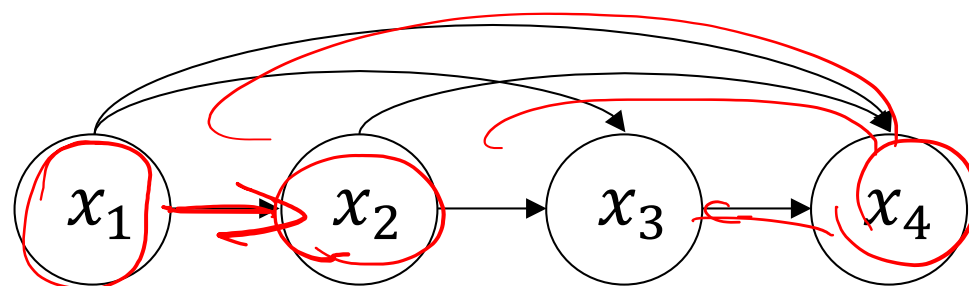


Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]

Resurgence of deep generative models

- Variational autoencoders (VAEs) [Kingma & Welling, 2014]
/ Neural Variational Inference and Learning (NVIL) [Mnih & Gregor, 2014]
- Generative adversarial networks (GANs) [Goodfellow et al., 2014]
- Generative moment matching networks (GMMNs) [Li et al., 2015; Dziugaite et al., 2015]
- Autoregressive neural networks
- Reversible architectures
- Diffusion models

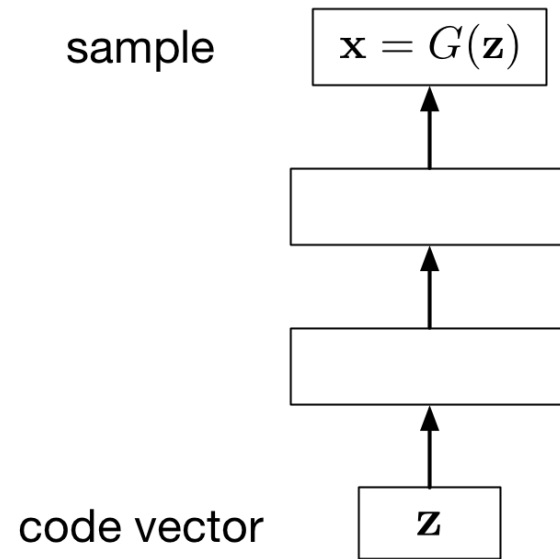


Generative Adversarial Networks

Recap: Implicit Generative Models

$x = G(z)$

- **Implicit generative models** implicitly define a probability distribution
- Start by sampling the **code vector** \mathbf{z} from a fixed, simple distribution (e.g. spherical Gaussian)
- The **generator network** computes a differentiable function G mapping \mathbf{z} to an \mathbf{x} in data space

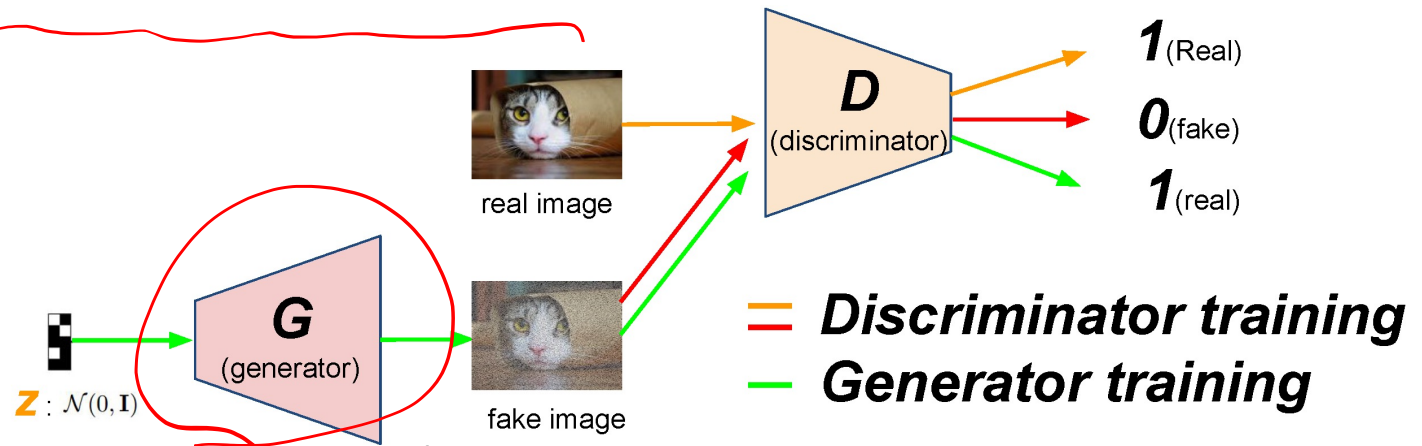


- a stochastic process to simulate data \mathbf{x}
- Intractable to evaluate likelihood

Recap: Generative Adversarial Nets (GANs)

- Learning
 - A **minimax** game between the generator and the discriminator
 - Train D to maximize the probability of assigning the correct label to both training examples and generated samples
 - Train G to fool the discriminator

$$\max_D \mathcal{L}_D = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))]$$
$$\min_G \mathcal{L}_G = \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))].$$



Recap: Optimality of GANs

- Objectives:

$$\max_D \mathcal{L}_D = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))]$$

$$\min_G \mathcal{L}_G = \mathbb{E}_{\mathbf{x} \sim G(\mathbf{z}), \mathbf{z} \sim p(\mathbf{z})} [\log(1 - D(\mathbf{x}))].$$

- Global optimality: $p_g = p_{data}$

Handwritten notes in red ink:

KL

$KL(p_{\theta}, p_{data})$

KL

Recap: A better loss function: non-saturating GAN

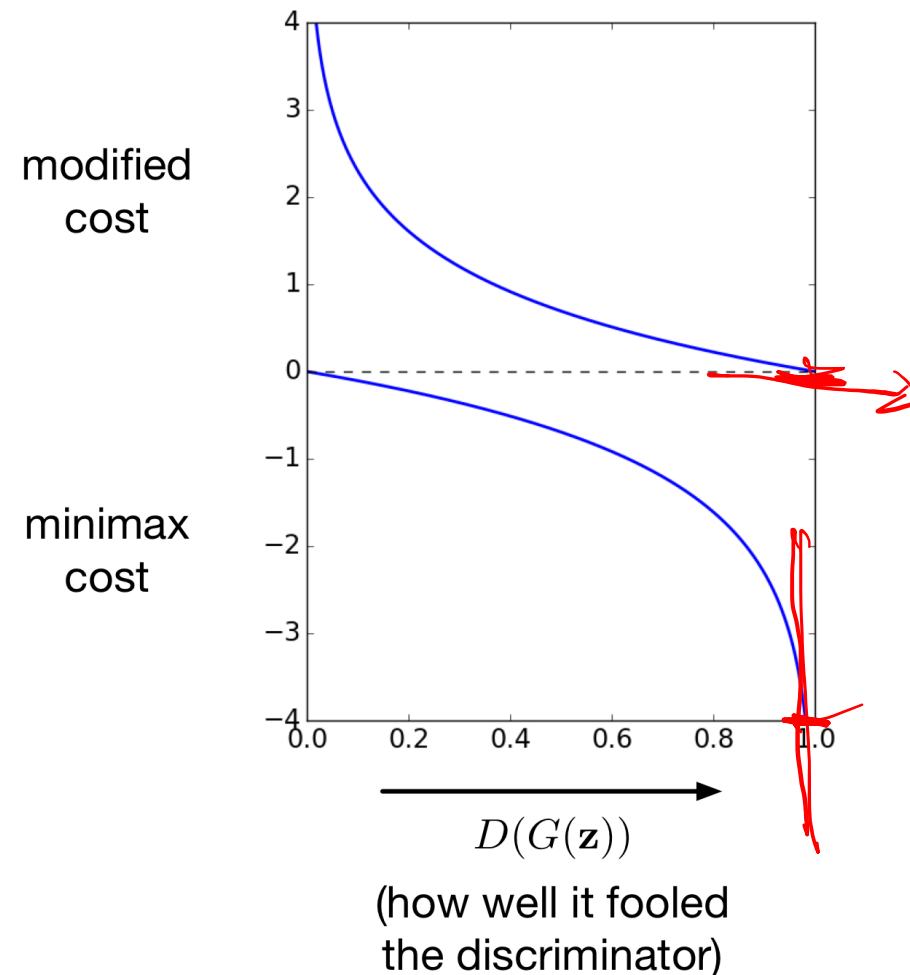
- Original minimax cost:

$$\mathcal{J}_G = \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$$

- Modified generator cost:

$$\mathcal{J}_G = \mathbb{E}_{\mathbf{z}}[-\log D(G(\mathbf{z}))]$$

- This fixes the saturation problem.



Wasserstein GAN (WGAN)

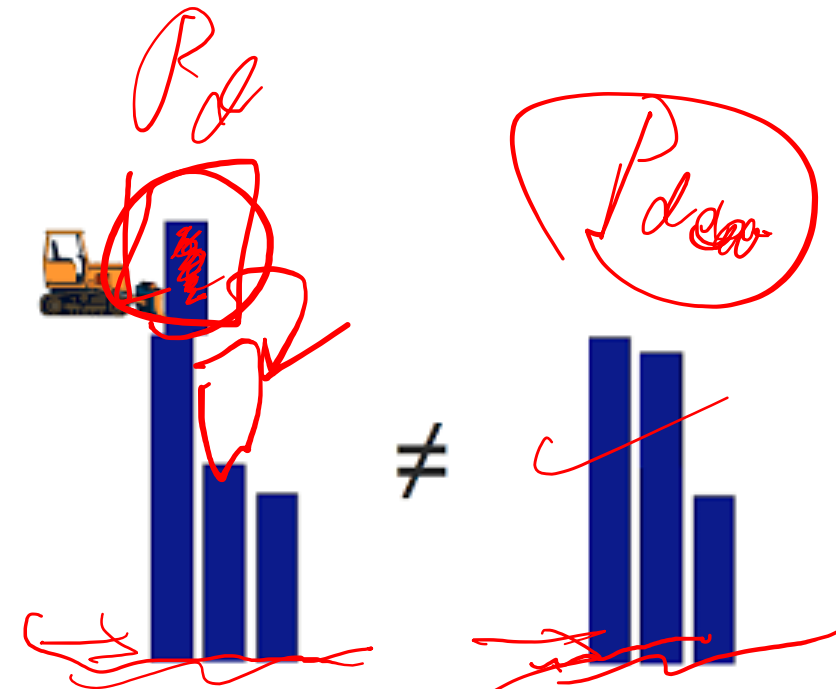
- If our data are on a **low-dimensional** manifold of a high dimensional space, the model's manifold and the true data manifold can have a **negligible intersection in practice**

Wasserstein GAN (WGAN)

- If our data are on a **low-dimensional** manifold of a high dimensional space, the model's manifold and the true data manifold can have a **negligible intersection in practice**
- The loss function and gradients may not be continuous and well behaved

Wasserstein GAN (WGAN)

- If our data are on a **low-dimensional** manifold of a high dimensional space, the model's manifold and the true data manifold can have a **negligible intersection in practice**
- The loss function and gradients may not be continuous and well behaved
- The **Wasserstein Distance** is well defined
 - Earth Mover's Distance
 - Minimum transportation cost for making one pile of dirt in the shape of one probability distribution to the shape of the other distribution



Wasserstein GAN (WGAN)

- Objective

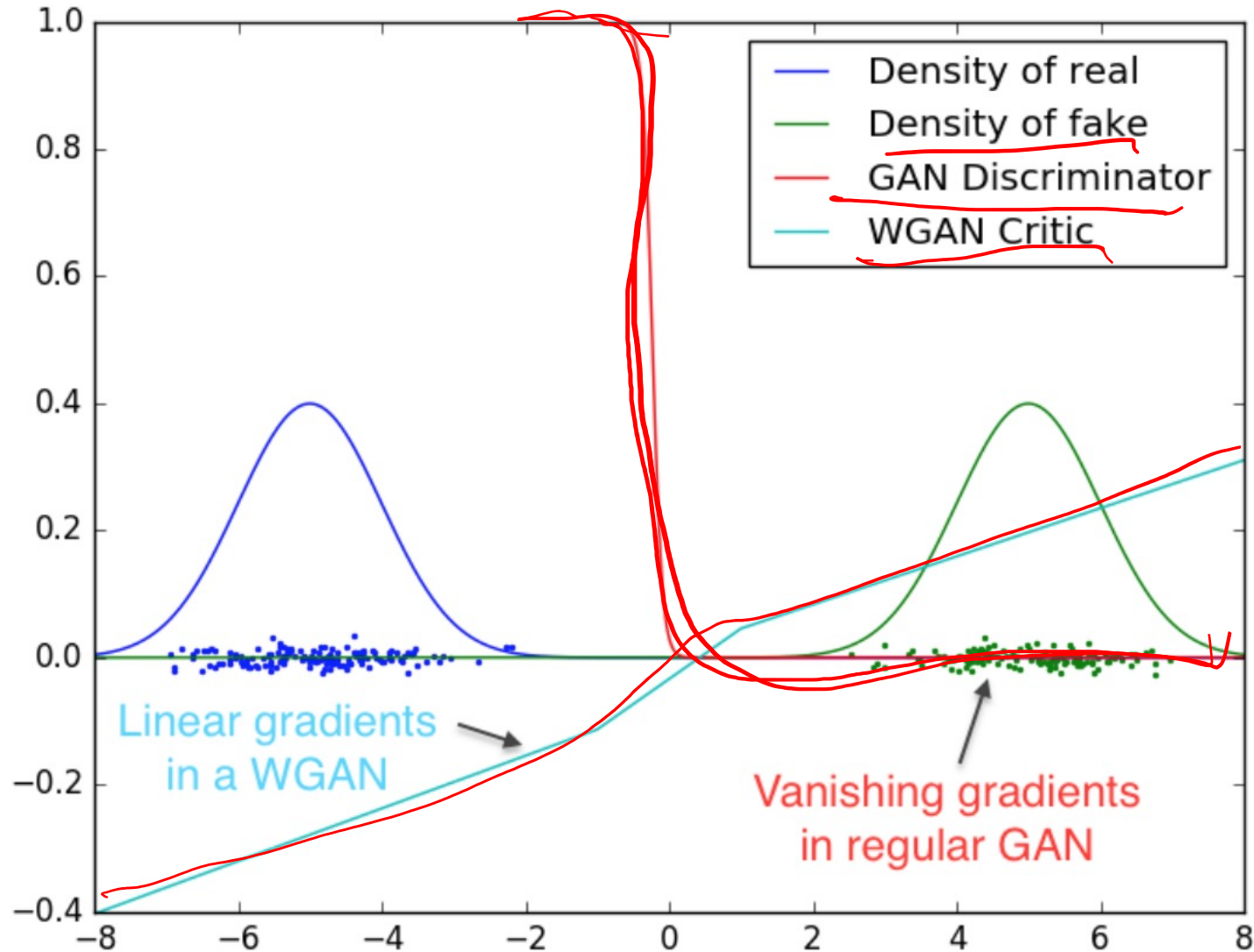
min

p_g

$$W(p_{data}, p_g) = \frac{1}{K} \sup_{\|D\|_L \leq K} \mathbb{E}_{x \sim p_{data}} [D(x)] - \mathbb{E}_{x \sim p_g} [D(x)]$$

- $\|D\|_L \leq K$: K - Lipschitz continuous
- Use gradient-clipping to ensure D has the Lipschitz continuity

WGAN vs Vanilla GAN



Standard Equation and GANs

- Recall SE:

$$\min_{q, \theta} -\alpha \mathbb{H}(q) + \beta \mathbb{D} \left(q(\mathbf{x}), p_{\theta}(\mathbf{x}) \right) - \mathbb{E}_{q(\mathbf{x})} \left[f(\mathbf{x}) \right]$$

- In MLE, f is a fixed function

$$f := f_{data}(\mathbf{x}; \mathcal{D}) = \log \mathbb{E}_{\mathbf{x}^* \sim \mathcal{D}} [\mathbb{1}_{\mathbf{x}^*}(\mathbf{x})]$$

- Intuitively, see f as a similarity metric that measures similarity of sample \mathbf{x} against real data \mathcal{D}
- Instead of the above manually fixed metric, can we learn a metric f_{ϕ} ?

Standard Equation and GANs

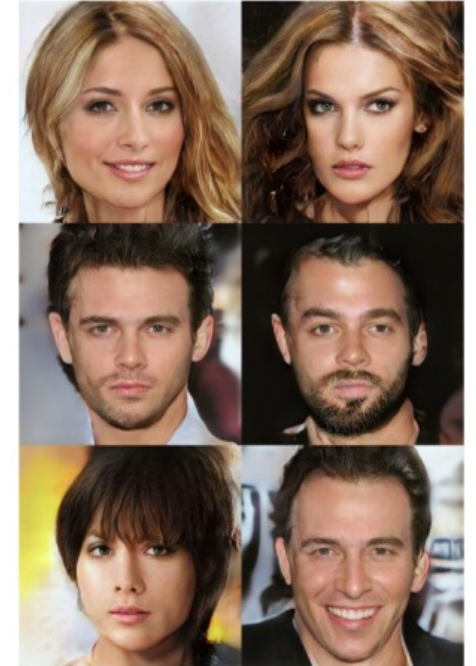
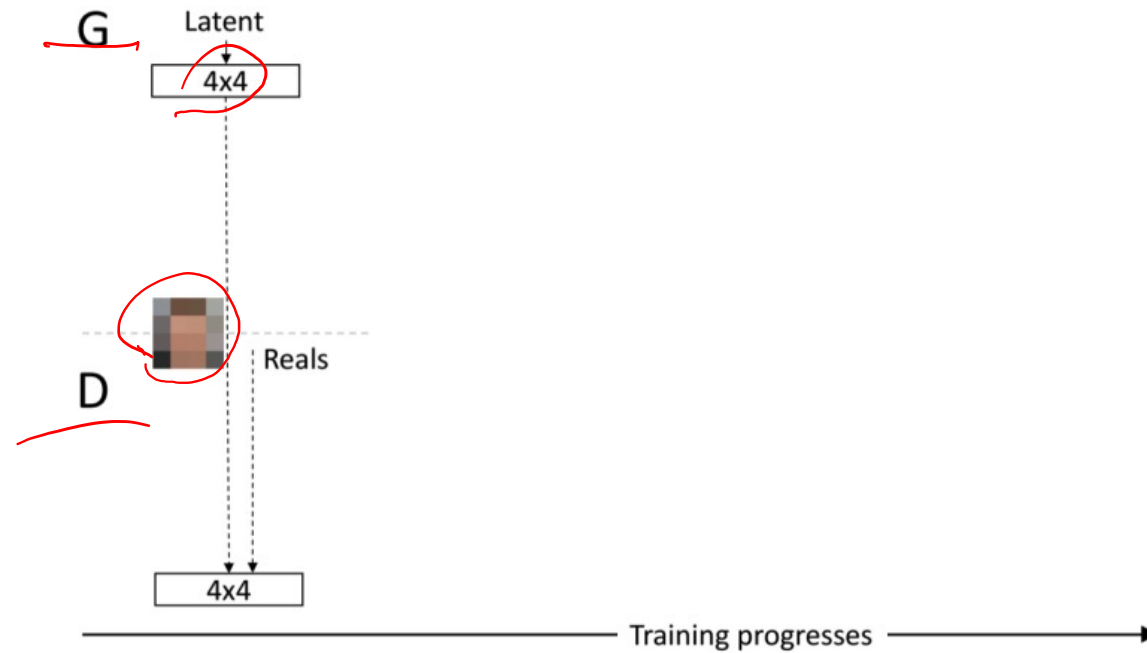
- Augment the standard objective to account for ϕ :

$$\min_{\theta} \max_{\phi} \min_q -\alpha \mathbb{H}(q) + \beta \mathbb{D} \left(q(\mathbf{x}), p_{\theta}(\mathbf{x}) \right) - \mathbb{E}_{q(\mathbf{x})} \left[f_{\phi}(\mathbf{x}) \right] + \mathbb{E}_{p_d(\mathbf{x})} \left[f_{\phi}(\mathbf{x}) \right]$$

- Set $\alpha = 0, \beta = 1$. Under mild conditions, the objective recovers:
 - Vanilla GAN [Goodfellow et al., 2014], when \mathbb{D} is JS-divergence and f_{ϕ} is a binary classifier
 - f -GAN [Nowozin et al., 2016], when \mathbb{D} is f -divergence
 - W-GAN [Arjovsky et al., 2017], when \mathbb{D} is Wasserstein distance and f_{ϕ} is a 1-Lipschitz function

Progressive GAN

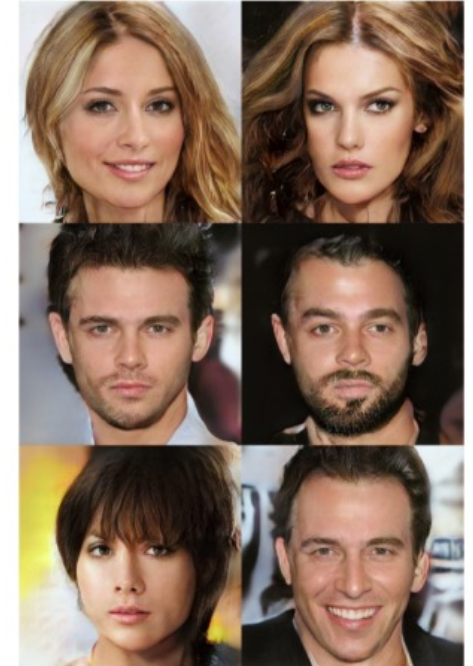
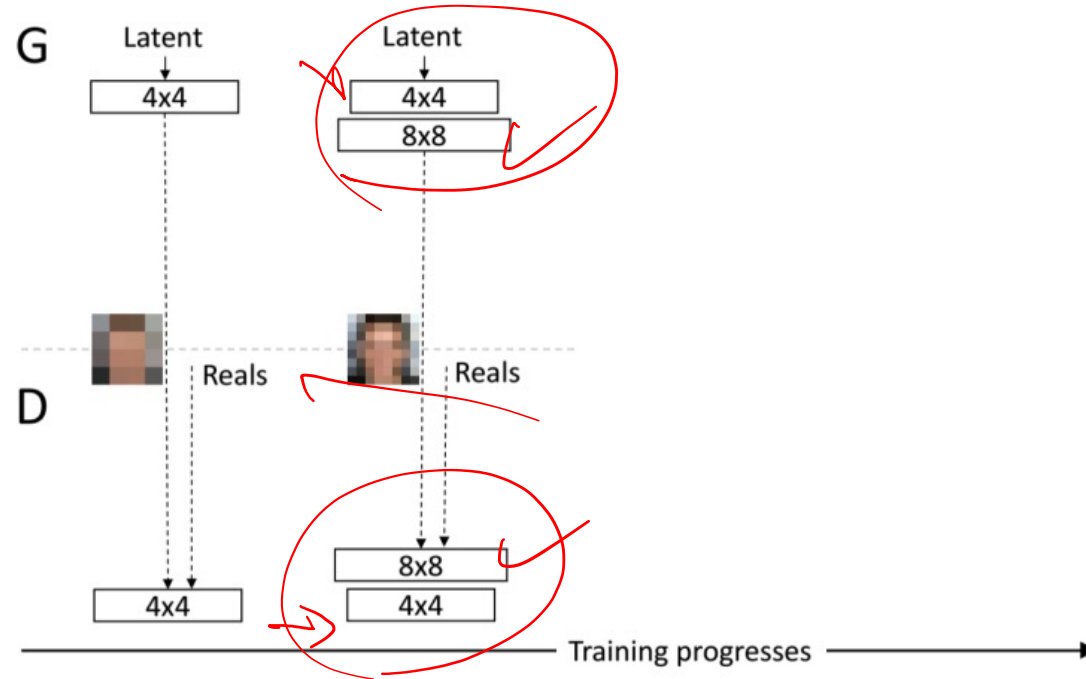
Low resolution images



Progressive GAN

Low resolution images

add in
additional
layers



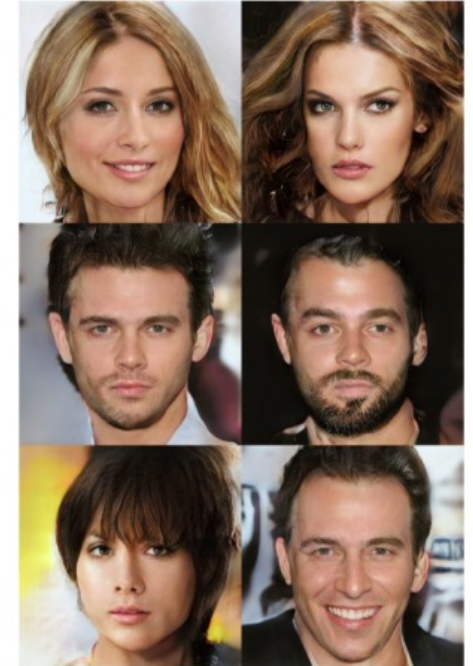
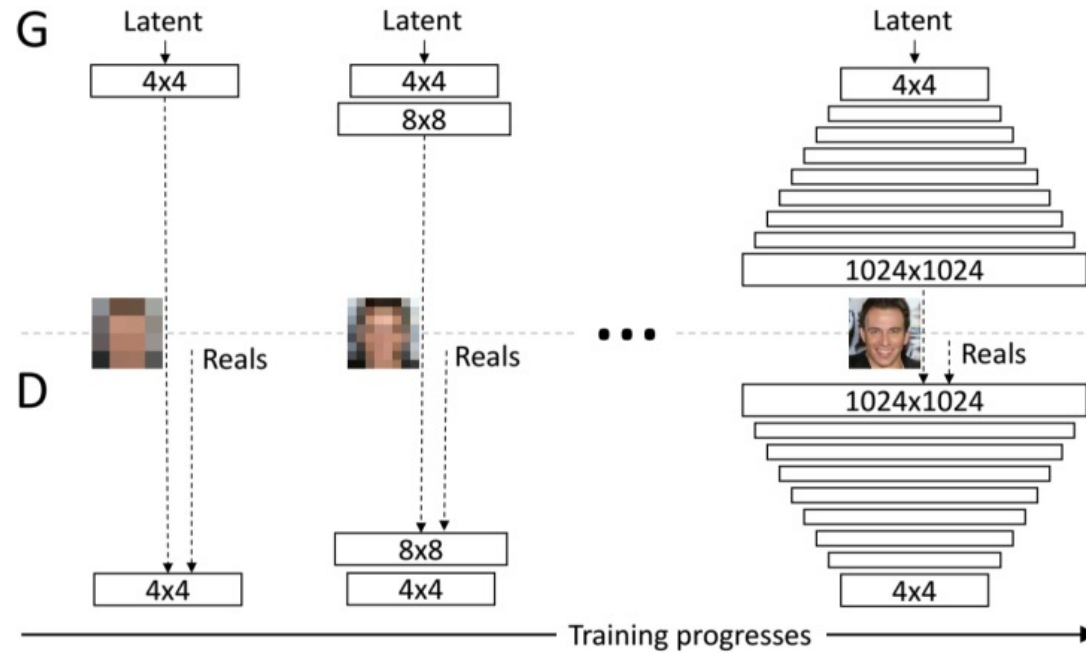
Progressive GAN

Low resolution images

add in
additional
layers



High resolution images



BigGAN

BigGAN

- GANs benefit dramatically from **scaling**

BigGAN

- GANs benefit dramatically from **scaling**
- 2x – 4x more parameters
- 8x larger batch size
- Simple architecture changes that improve scalability

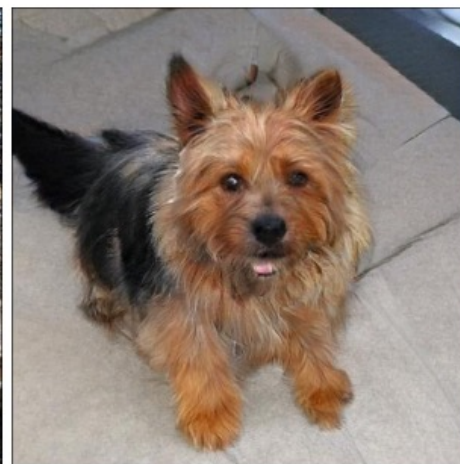
BigGAN

- GANs benefit dramatically from **scaling**
- 2x – 4x more parameters
- 8x larger batch size
- Simple architecture changes that improve scalability



BigGAN

- GANs benefit dramatically from **scaling**
- 2x → 4x more parameters
- 8x → 16x more parameters
- Similar quality to human



Normalizing Flow (NF)

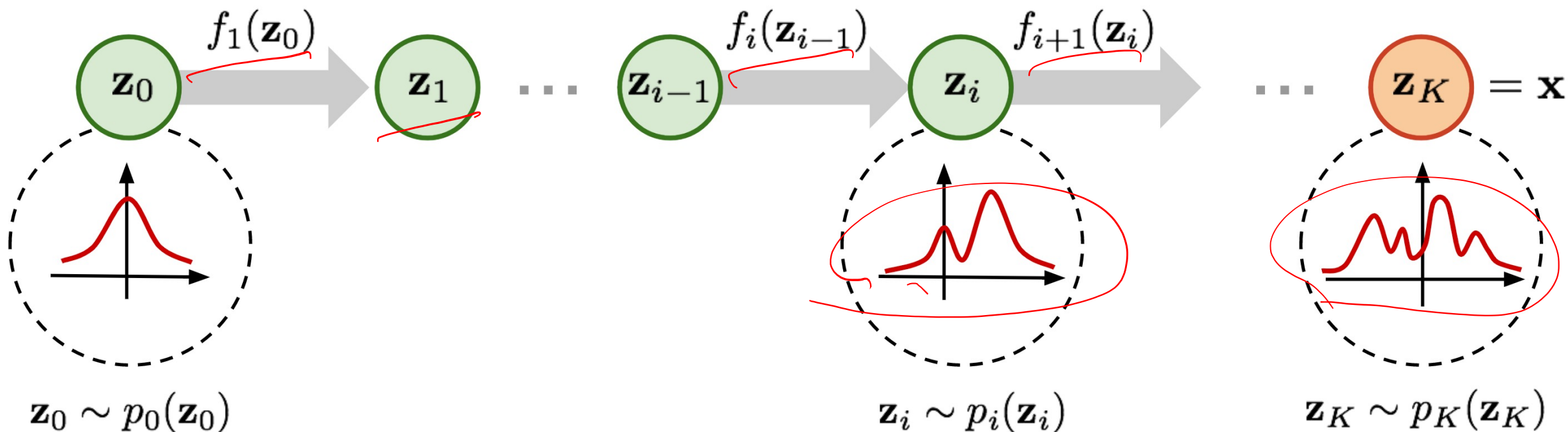
Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of transformation functions

$p(x)$ \rightarrow $p(z)$

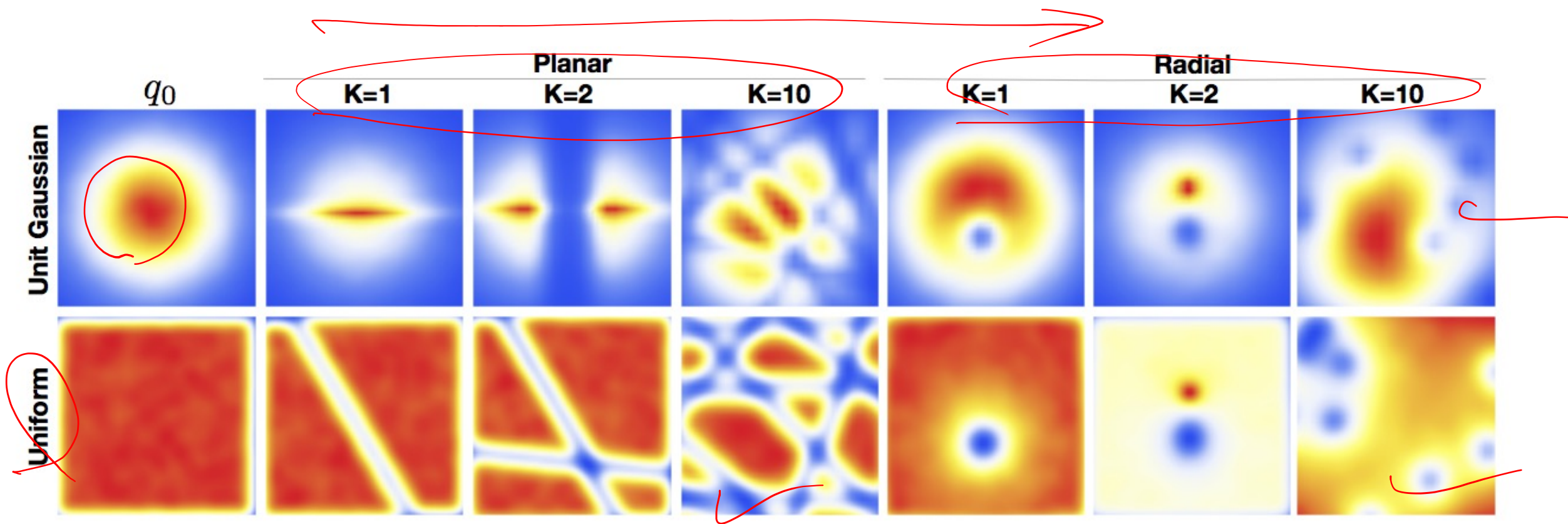
Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**



Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**



Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

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

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

Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

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

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

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

Transformation function f

- > • Invertible

Normalizing Flow (NF)



- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

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

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

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

density: $p(\mathbf{x}) = p(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right|$
 $= p(f^{-1}(\mathbf{x})) \left| \det \frac{df^{-1}}{d\mathbf{x}} \right|$

$$\det \frac{df^{-1}}{d\mathbf{x}} \text{ -- Jacobian determinant}$$

Transformation function f

-----> • Invertible

Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

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

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

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

density:
$$p(\mathbf{x}) = p(\mathbf{z}) \left| \det \frac{d\mathbf{z}}{d\mathbf{x}} \right|$$
$$= p(f^{-1}(\mathbf{x})) \left| \det \frac{df^{-1}}{d\mathbf{x}} \right|$$

Transformation function f

-----> • Invertible

-----> • Jacobian determinant easy to compute
e.g., choose $df^{-1}/d\mathbf{x}$ to be a triangular matrix

$$\det \frac{df^{-1}}{d\mathbf{x}} \text{ -- Jacobian determinant}$$

Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

$$\mathbf{z}_0 \sim p(\mathbf{z}_0)$$

$$\mathbf{x} = \mathbf{z}_K = f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$$

Transformation function f_i

inference: $\mathbf{z}_i = f_i^{-1}(\mathbf{z}_{i-1})$

density: $p(\mathbf{z}_i) = p(\mathbf{z}_{i-1}) \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$

-----> • Invertible

-----> • Jacobian determinant easy to compute
e.g., choose $df_i^{-1}/d\mathbf{z}_i$ to be a triangular matrix

Normalizing Flow (NF)

- Transforms a simple distribution into a complex one by applying a sequence of **transformation functions**

$$\mathbf{z}_0 \sim p(\mathbf{z}_0)$$

$$\mathbf{x} = \mathbf{z}_K = f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$$

Transformation function f_i

inference: $\mathbf{z}_i = f_i^{-1}(\mathbf{z}_{i-1})$

----->

- Invertible

density: $p(\mathbf{z}_i) = p(\mathbf{z}_{i-1}) \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$

----->

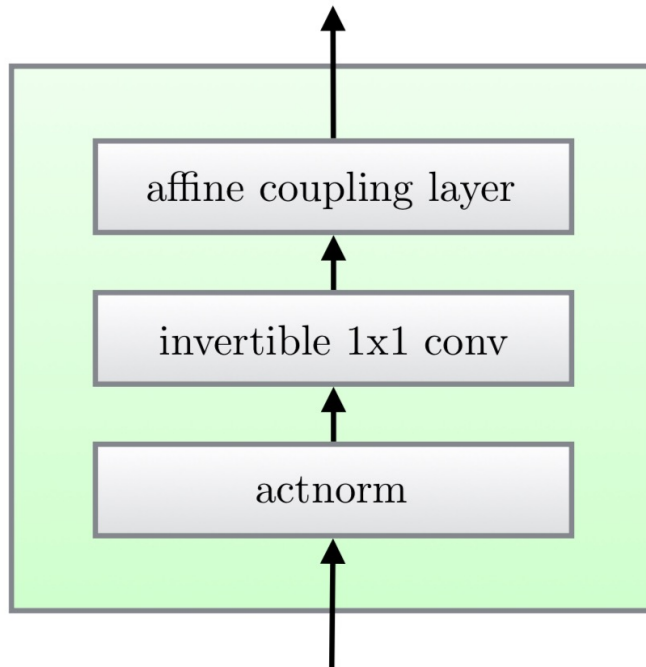
- Jacobian determinant easy to compute
e.g., choose $df_i^{-1}/d\mathbf{z}_i$ to be a triangular matrix

training: maximizes data log-likelihood

$$\log p(\mathbf{x}) = \log p(\mathbf{z}_0) + \sum_{i=1}^K \log \left| \det \frac{d\mathbf{z}_{i-1}}{d\mathbf{z}_i} \right|$$

GLOW

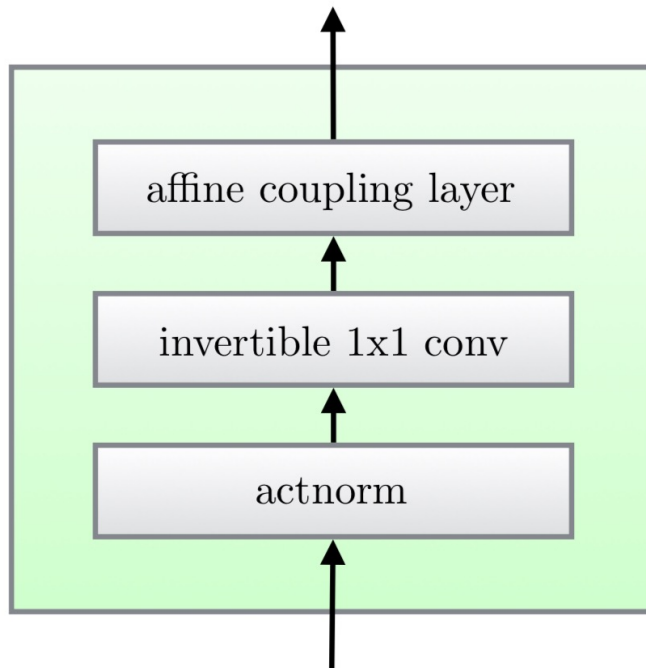
- [Kingma and Dhariwal., 2018]



One step of flow in the Glow model

GLOW

- [Kingma and Dhariwal., 2018]



One step of flow in the Glow model

Key Takeaways

- GANs:
 - Implicit generative model
 - Minimax formulation
 - non-saturating GANs
 - WGAN
- Normalizing Flow
 - Transforms a simple distribution into a complex one by applying a sequence of transformation functions

Questions?