

An Introduction to Deep Generative Models

Sajjad Amini

Manning College of Information and Computer Sciences
University of Massachusetts Amherst

November 7, 2023

Contents

- 1 The Quest for Deep Generative Models
- 2 Big Picture
- 3 Generative Adversarial Networks
(Latent Variable Model)
 - Family Selection
 - Distance Metric
 - Applications
- 4 Autoregressive Models
(Likelihood Based Model)
 - Family Selection
 - Distance Metric
 - Sampling
 - Application
- 5 Latent Variable Models with Likelihood Training
- 6 Summary

Section 1

The Quest for Deep Generative Models

IRIS dataset

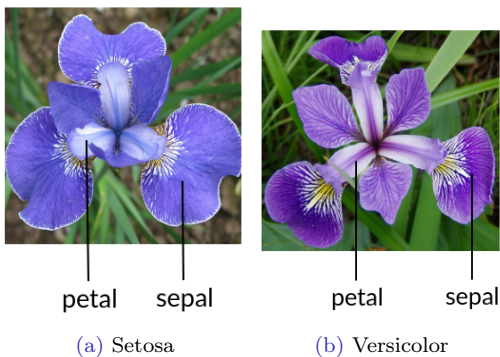


Figure: Two types of Iris flower

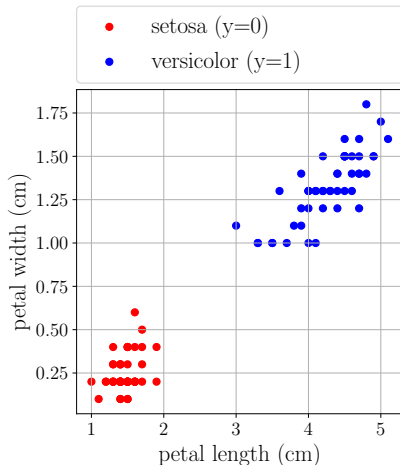


Figure: Dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

Simple Generative Model: Binary Logistic Regression

Model

Model:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y | \overbrace{\sigma(\mathbf{w}^T \mathbf{x} + b)}^p)$$
$$\Rightarrow \begin{cases} p_{\theta}(y=1|\mathbf{x}) = p \\ p_{\theta}(y=0|\mathbf{x}) = 1-p \end{cases}$$

where:

$\sigma(\cdot)$	Sigmoid function
\mathbf{w}	Weight vector
b	Bias value
$\theta = [b; \mathbf{w}]$	Model parameters

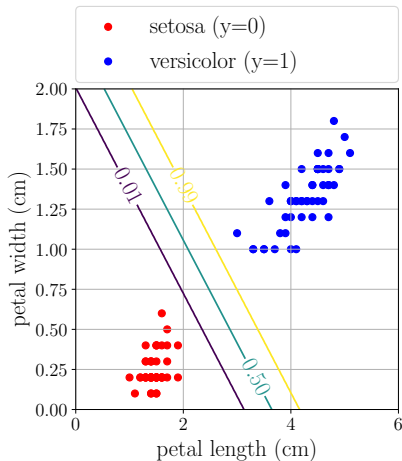


Figure: BLR contours for p

Sampling from Binary Logistic Regression

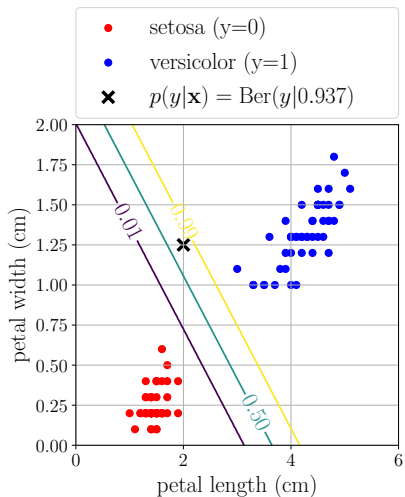


Figure: Query inference

Sampling

Sampling from $p_{\theta}(y|\mathbf{x}_q) = \text{Ber}(y|0.937)$

- $u \sim \mathcal{U}(0, 1)$
- $y = \begin{cases} 0 & \text{if } u \geq 0.937 \\ 1 & \text{if } u < 0.937 \end{cases}$

Limitation

👉 $y \in \{0, 1\}$: y is a binary random variable.

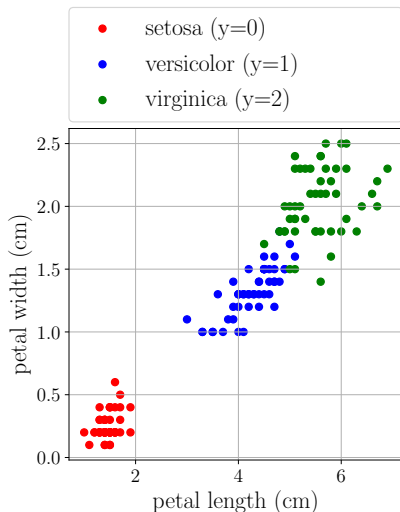


Figure: Dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$

Generative Model: Multinomial Logistic Regression

Model

Model:

$$p_{\theta}(y|\mathbf{x}) = \text{Cat}(y | \underbrace{\mathcal{S}(\mathbf{W}\mathbf{x} + \mathbf{b})}_{\mathbf{p}=[p_0, \dots, p_{L-1}]^T})$$
$$\Rightarrow \begin{cases} p_{\theta}(y = 0|\mathbf{x}) = p_0 \\ \vdots \\ p_{\theta}(y = L - 1|\mathbf{x}) = p_{L-1} \end{cases}$$

where:

$\mathcal{S}(\cdot)$	Softmax function
$\mathbf{W} \in \mathbb{R}^{C \times D}$	Weight matrix
$\mathbf{b} \in \mathbb{R}^C$	Bias vector
$\theta = \{\mathbf{W}, \mathbf{b}\}$	Model parameters
$\mathbf{a} = \mathbf{W}\mathbf{x} + \mathbf{b}$	logits vector

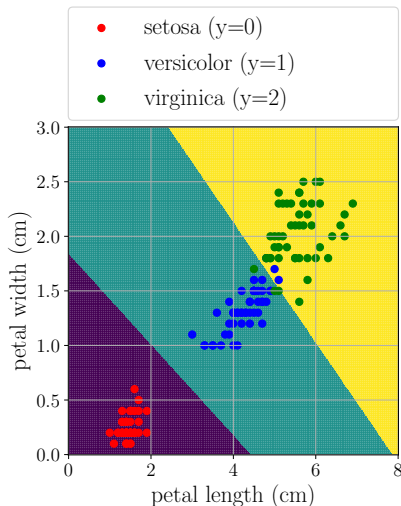
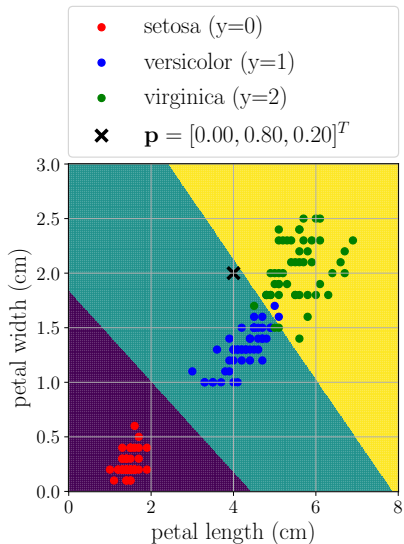


Figure: BLR contours for p

Sampling from Multinomial Logistic Regression



Sampling

Sampling from $p_{\theta}(y|\mathbf{x}_q)$ =

Cat $\left(y \mid \begin{bmatrix} 0.00 \\ 0.80 \\ 0.20 \end{bmatrix} \right)$

$$\bullet \mathbf{p} = \begin{bmatrix} 0.0 \\ 0.2 \\ 0.8 \end{bmatrix} \Rightarrow \mathbf{c} = \begin{bmatrix} 0.0 \\ 0.2 \\ 1 \end{bmatrix}$$

$$\bullet u \sim \mathcal{U}(0, 1)$$

$$\bullet y = \begin{cases} 0 & \text{if } u \leq 0 \\ 1 & \text{if } 0 < u \leq 0.2 \\ 2 & \text{if } u > 0.2 \end{cases}$$

Figure: Query inference

Limitation

Linear decision boundaries

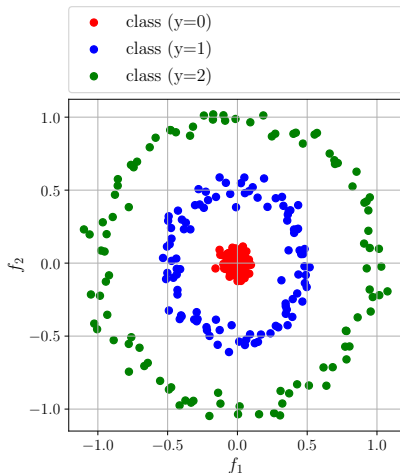
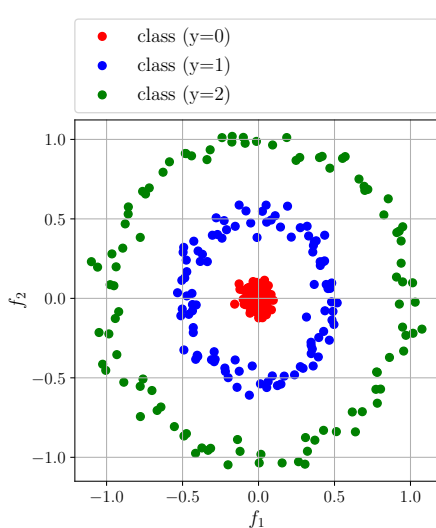
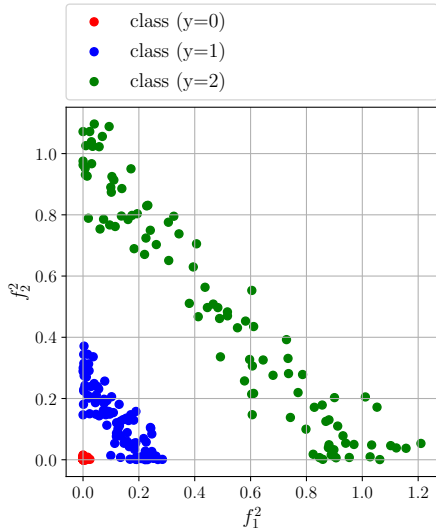


Figure: Not linearly separable dataset

Feature Transformation

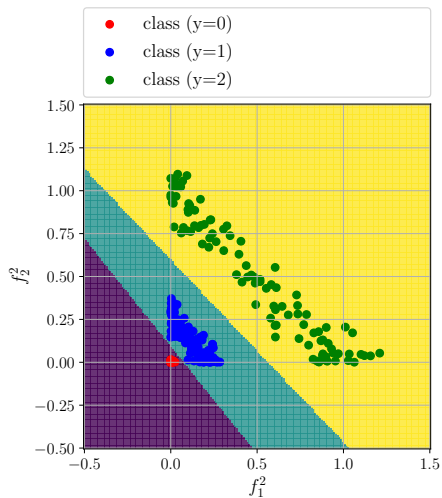


(a) Original features

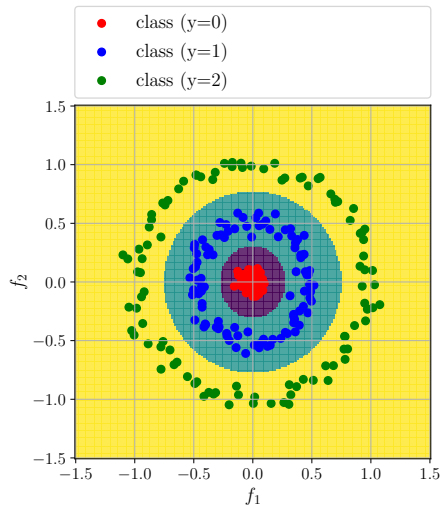


(b) Transformed features

Classification in the Transformed Space



(a) MLR result in the transformed space



(b) Transformed features

The Era of Deep Learning

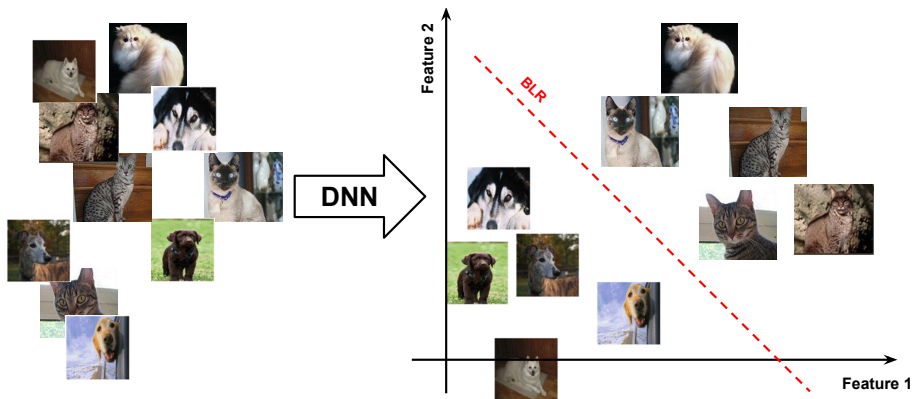


Figure: Learning suitable transformation from data using Deep Neural Networks

What We Can Do So Far

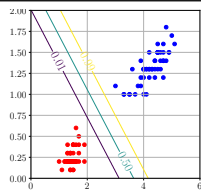
Model

Distribution

Sample

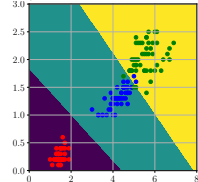
Binary Logistic Regression

$$p_{\theta}(y|\mathbf{x})$$
$$\begin{cases} \mathbf{x} \in \mathbb{R}^D \\ y \in \{0, 1\} \end{cases}$$



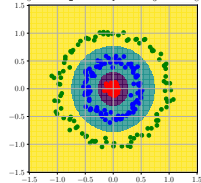
Multinomial Logistic Regression

$$p_{\theta}(y|\mathbf{x})$$
$$\begin{cases} \mathbf{x} \in \mathbb{R}^D \\ y \in \{1, 2, \dots, L\} \end{cases}$$



DNN Feature Transformation
+
Multinomial Logistic Regression

$$p_{\theta}(y|\mathbf{x})$$
$$\begin{cases} \mathbf{x} \in \mathbb{R}^D \\ y \in \{1, 2, \dots, L\} \end{cases}$$



Limitation

y is assumed to be a one-dimensional random variable and cannot be extended to the case where \mathbf{y} is a high-dimensional random vector.

Text-to-Speech Models

$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{An audio file} \\ \mathbf{x} : \text{A text} \end{cases}$$

$\mathbf{x} =$ “A single Wavenet can capture the characteristics of many different speakers with equal fidelity, not it’s fast.” $\xrightarrow{\text{Sampling } p(\mathbf{x}|\mathbf{y})}$ $\mathbf{y} =$

Text-to-Image Models

$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{An image} \\ \mathbf{x} : \text{A text} \end{cases}$$



Figure: \mathbf{x} for \mathbf{y} = “Teddy bears swimming at the Olympics 400m Butterfly event.”

Image Colorization

$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{A Colored image} \\ \mathbf{x} : \text{A Gray - scale image} \end{cases}$$



(a) \mathbf{x}



(b) \mathbf{y}



(c) Ground truth

Image Inpainting

$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{A clean image} \\ \mathbf{x} : \text{A corrupted image} \end{cases}$$



(a) \mathbf{x}



(b) \mathbf{y}



(c) Ground truth

Image Uncropping

$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{A clean image} \\ \mathbf{x} : \text{A cropped image} \end{cases}$$



(a) \mathbf{x}



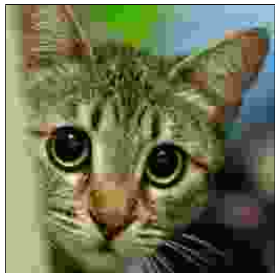
(b) \mathbf{y}



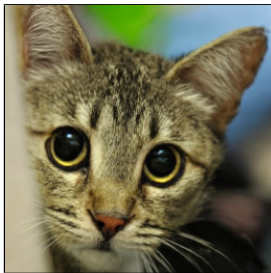
(c) Ground truth

Image Restoration

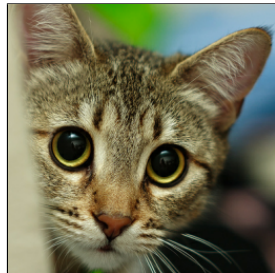
$$p(\mathbf{y}|\mathbf{x}) : \begin{cases} \mathbf{y} : \text{A clean image} \\ \mathbf{x} : \text{A degraded image} \end{cases}$$



(a) \mathbf{x}



(b) \mathbf{y}



(c) Ground truth

The Era of Deep Learning

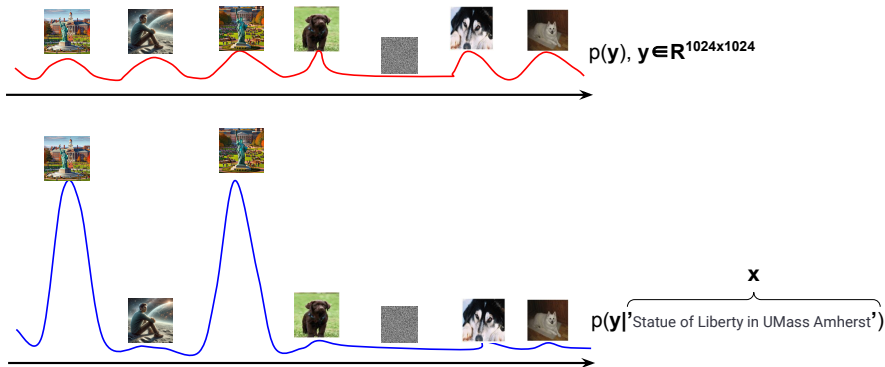
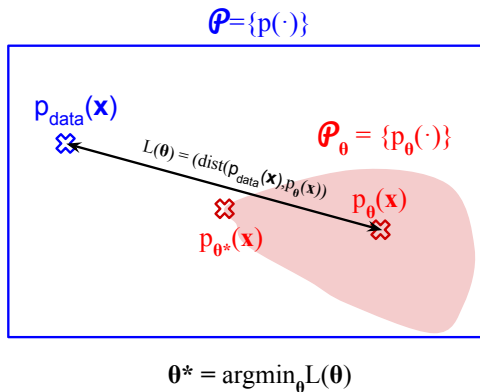
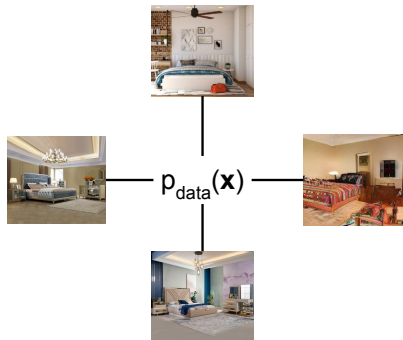


Figure: Probabilistic intuition to deep generative models

Section 2

Big Picture

Deep Generative Modeling Big Picture [4]



Generative Model Learning

- **Experience:**

$$\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$$

- **Objective:**

- Estimating distribution $p_\theta(\mathbf{x})$

- **Tasks:**

- Density estimation: $p_\theta(\mathbf{x}_{\text{new}})$
- Generation: $\mathbf{x}_{\text{new}} \sim p_\theta(\mathbf{x})$
- High-level representation \mathbf{z}

Section 3

Generative Adversarial Networks (Latent Variable Model)

Subsection 1

Family Selection

Family Overview

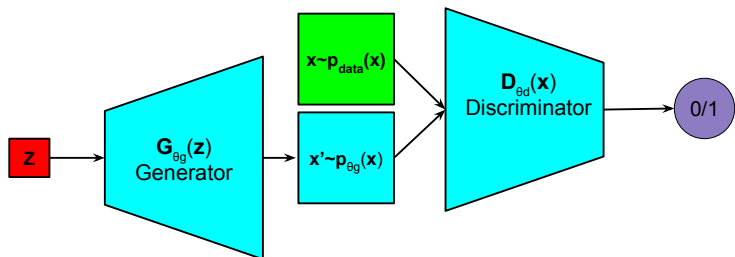


Figure: GAN Architecture

Specifications

- $z \in \mathbb{R}^{100} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- $x' \in \mathbb{R}^{64 \times 64 \times 3}$: Generated Images
- $x \in \mathbb{R}^{64 \times 64 \times 3}$: Real Images
- θ_g, θ_d : Generator and discriminator parameters

Generator in Deep Convolutional GAN [5]

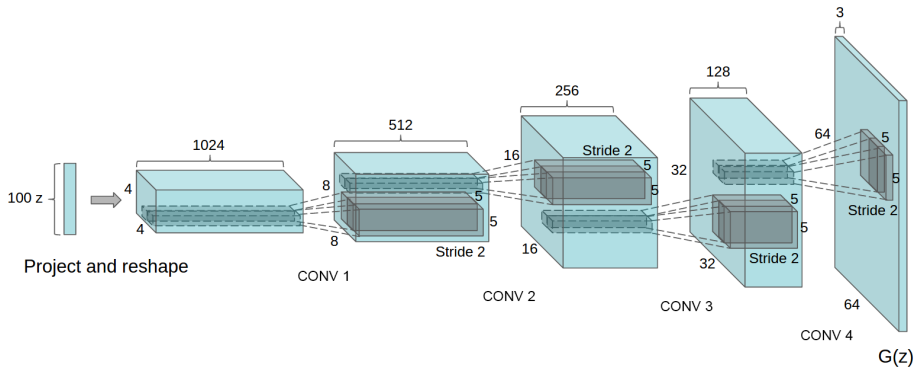


Figure: Generator Architecture in DCGAN (input: $z \in \mathbb{R}^{100}$, output: $x' \in \mathbb{R}^{64 \times 64 \times 3}$)

Discriminator in Deep Convolutional GAN [5]

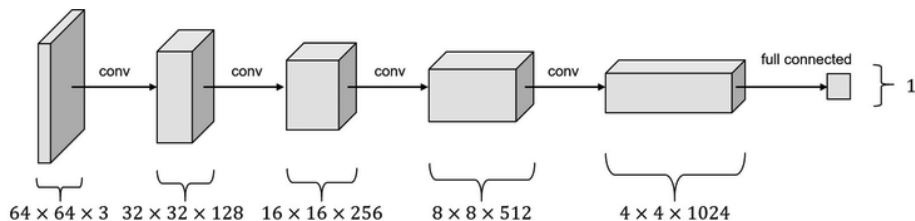


Figure: Discriminator Architecture in DCGAN (input: $\mathbf{x} \in \mathbb{R}^{64 \times 64 \times 3}$, output: $p \in [0, 1]$) [6]

Subsection 2

Distance Metric

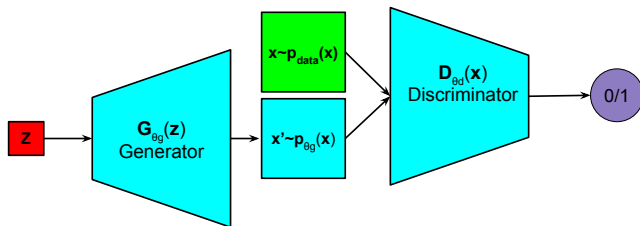


Figure: GAN Architecture

Training Generator

Assume we have trained a powerful discriminator such that:

$$\begin{cases} \mathcal{D}_{\theta_d}(\mathbf{x}) \simeq 1 \\ \mathcal{D}_{\theta_d}(\mathbf{x}') \simeq 0 \end{cases}$$

Then a good generator can be trained as:

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(\overbrace{G_{\theta_g}(z)}^{x'}))$$

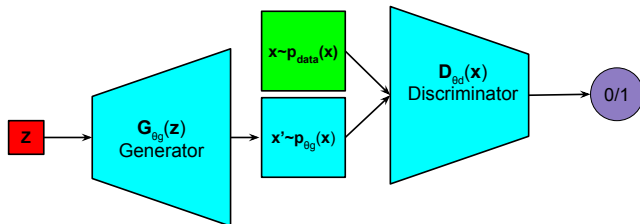


Figure: GAN Architecture

Training Discriminator

Assume we have trained a powerful generator that can generate samples quite similar to real ones. Then a good discriminator can be trained as:

$$\max_{\theta_d} \left[\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \log (1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))) \right]$$

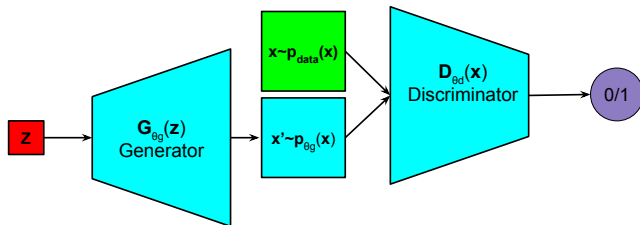


Figure: GAN Architecture

Training GAN

In practice, there is neither a powerful generator nor a powerful discriminator available at the start. Thus both should be trained in a min-max optimization as:

$$\operatorname{argmin}_{\theta_g} \operatorname{argmax}_{\theta_d} \left[\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D_{\theta_d}(\mathbf{x}) + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \log (1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}))) \right]$$

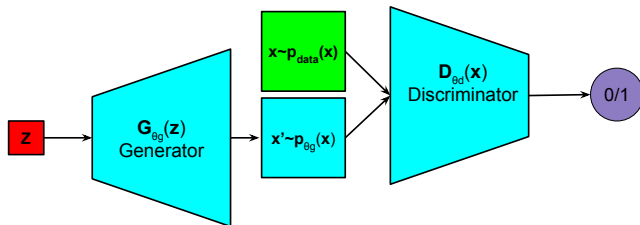


Figure: GAN Architecture

Desired Convergence situation

We have a successful learning if:

- A powerful discriminator is trained.
- The generator can fool the discriminator with high probability.

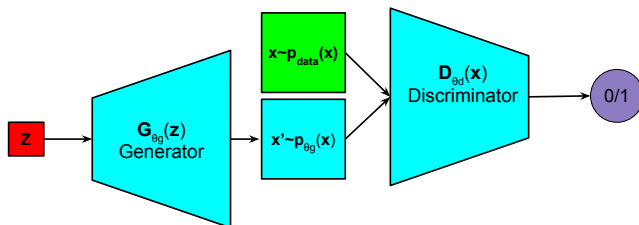


Figure: GAN Architecture

Monte Carlo Estimation

$$\operatorname{argmin}_{\theta_g} \operatorname{argmax}_{\theta_d} \left[\underbrace{\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D_{\theta_d}(\mathbf{x})}_{\frac{1}{K} \sum_{k=1}^K \log D_{\theta_d}(\mathbf{x}_k)} + \underbrace{\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} \log (1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z})))}_{\frac{1}{K} \sum_{k=1}^K \log (1 - D_{\theta_d}(G_{\theta_g}(\mathbf{z}_k)))} \right]$$

Subsection 3

Applications

Samples from Deep Convolutional GAN



Figure: Generated sample images from DCGAN trained on bedroom images [5]

Application of Latent Representation [5]

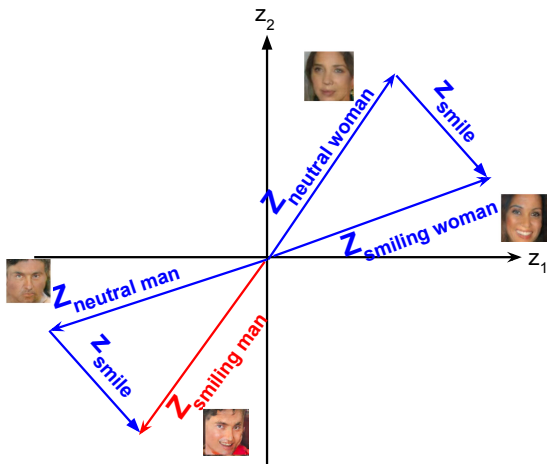


Figure: Finding the direction of *smile* in latent space z_{smile}

Application of Latent Representation

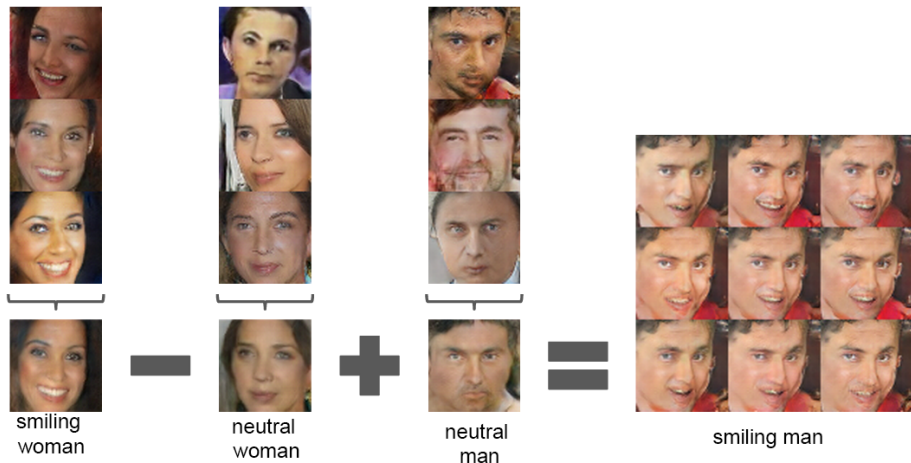


Figure: Finding the direction of *smile* in latent space z [5]

Application of Latent Representation

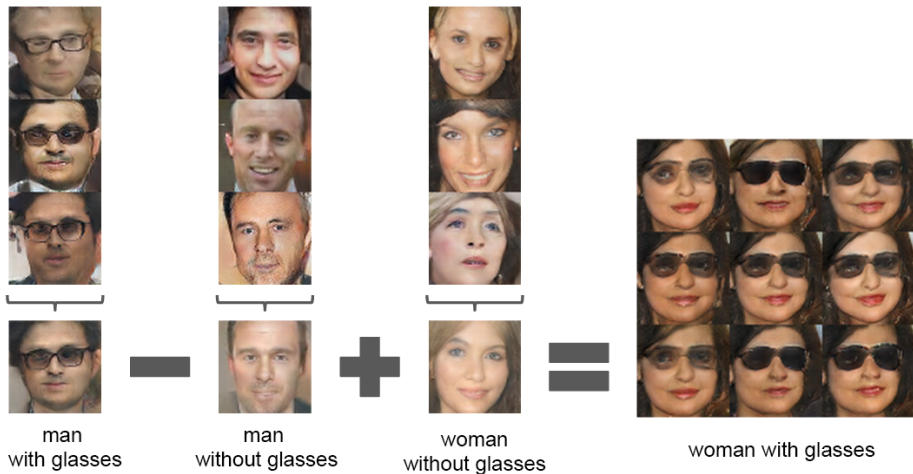


Figure: Finding the direction of *wearing glasses* in latent space z [5]

Section 4

Autoregressive Models (Likelihood Based Model)

Subsection 1

Family Selection

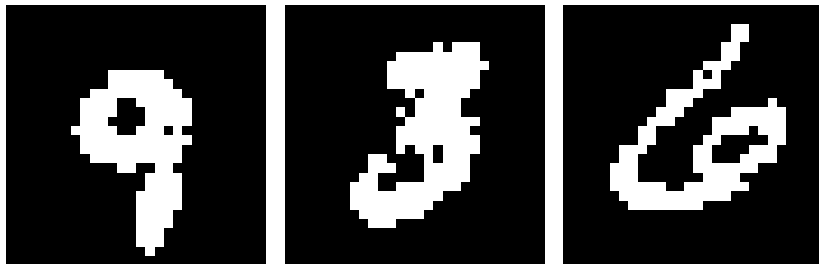
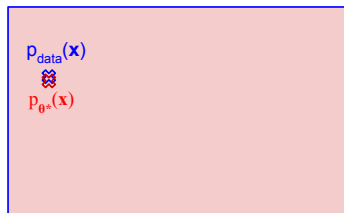


Figure: Samples from binary MNIST dataset ($\mathbf{x} \in \mathbb{R}^{784}$ and $x_i \in \{0, 1\}$)

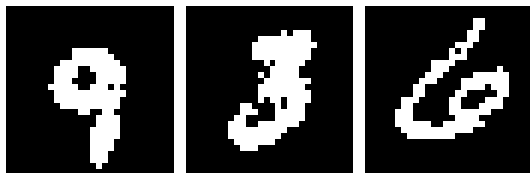
Full Joint Using Chain Rule

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1)p(x_3|x_2, x_1) \dots p(x_{784}|x_{783}, \dots, x_1)$$

Complete Distribution Search



(a) Typical \mathcal{P} and \mathcal{P}_θ



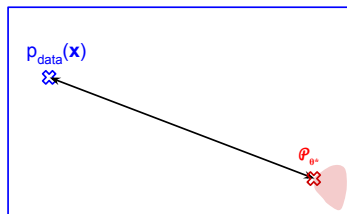
(b) Generated samples from p_{θ^*} assuming $\mathcal{P} = \mathcal{P}_\theta$ we can find p_{θ^*} efficiently

Challenges

$$p(\mathbf{x}) = \underbrace{p(x_1)}_{\#p=1} \underbrace{p(x_2|x_1)}_{\#p=2} \dots \underbrace{p(x_{784}|x_{783}, \dots, x_1)}_{\#p=2^{n-1}}$$

- Intractable number of parameters: $2^{784} - 1$
- Hard to optimize

Small Subset Search



(a) Typical \mathcal{P} and \mathcal{P}_θ



(b) Generated samples from p_{θ^*} assuming we can find p_{θ^*} efficiently

Challenges

$$p(\mathbf{x}) = \underbrace{p(x_1)}_{\#p=1} \underbrace{p(x_2)}_{\#p=1} \dots \underbrace{p(x_{784})}_{\#p=1}$$

- The resulting $p_\theta(\mathbf{x})$ cannot capture the data distribution.

Chain Rule

Using the chain rule, we have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_i|\mathbf{x}_{<i}) \dots p(x_{784}|\mathbf{x}_{<784}), \quad \mathbf{x}_{<i} \triangleq [x_1, \dots, x_{i-1}]^T$$

Resemblance to BLR

Consider a factor in the right-hand side of the chain rule as:

$$p(x_i|\mathbf{x}_{<i}), x_i \in \{0, 1\}$$

This problem can be easily solved via BLR (or BLR over DNN features) as:

$$p(x_i|\mathbf{x}_{<i}; \mathbf{w}_i, b_i) = \text{Ber}(x_i|\sigma(b_i + \mathbf{w}_i^T \mathbf{x}_{<i})), \quad \begin{cases} b_i \in \mathbb{R} \\ \mathbf{w}_i \in \mathbb{R}^{i-1} \end{cases}$$

Our New Parametric Family

Name

The new parametric distribution family is called **Deep Autoregressive Models**, as they calculate the current state x_i given its past $\mathbf{x}_{<i}$.

Number of Parameters

We have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_n|\mathbf{x}_{<n})$$
$$p(\mathbf{x}) = \underbrace{\text{Ber}(x_1|\sigma(b_1))}_{\#p=1} \underbrace{\text{Ber}(x_2|\sigma(b_2 + w_2^T \mathbf{x}_{<2}))}_{\#p=2} \dots \underbrace{\text{Ber}(x_n|\sigma(b_n + \mathbf{w}_n^T \mathbf{x}_{<n}))}_{\#p=n}$$

Thus

- #Parameters for \mathcal{P}_θ : $1 + 2 + 3 + \dots + n = \frac{n(n+1)}{2}$.
- $\theta = \{b_1, b_2, w_2, \dots, b_n, \mathbf{w}_n\}$

Subsection 2

Distance Metric

Model Likelihood Estimation

Model Likelihood Estimation

We are interested in solving the following problem:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x})]$$

Solution via Monte Carlo Estimate

Using Monte Carlo estimate we have:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\boldsymbol{\theta}}(\mathbf{x})] \simeq \frac{1}{|\mathcal{D}|} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

Thus:

$$\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \frac{1}{|\mathcal{D}|} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

Calculating $p_{\theta}(\mathbf{x}_i)$

Assume $\mathbf{x} \in \mathbb{R}^2$ and our model is:

$$p(\mathbf{x}) = \text{Ber}(x_1|\sigma(b_1)) \text{Ber}(x_2|\sigma(b_2 + w_2 \times x_1))$$

To compute the probability $\mathbf{x}_i = [1, 0]^T$ we have:

$$\begin{aligned} p(\mathbf{x}_i = \begin{bmatrix} 1 \\ 0 \end{bmatrix}) &= \text{Ber}(x_1 = 1|\sigma(b_1)) \text{Ber}(x_2 = 0|\sigma(b_2 + w_2 \times 1)) \\ &= \sigma(b_1) \times [1 - \sigma(b_2 + w_2 \times 1)] \end{aligned}$$

Subsection 3

Sampling

Sampling

Assume $\mathbf{x} \in \mathbb{R}^2$ and our trained model is (trained parameters are: $\theta^* = \{b_1^*, b_2^*, w_2^*\}$):

$$p(\mathbf{x}) = \text{Ber}(x_1 | \sigma(b_1^*)) \text{Ber}(x_2 | \sigma(b_2^* + w_2^* \times x_1))$$

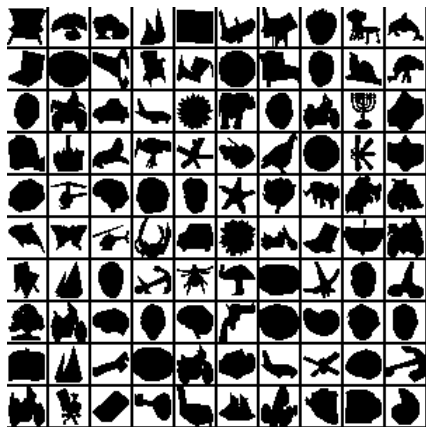
To sample we do:

- $\hat{x}_1 \leftarrow \text{Sample}\left(\text{Ber}(x_1 | \sigma(b_1^*))\right)$
- $\hat{x}_2 \leftarrow \text{Sample}\left(\text{Ber}(x_2 | \sigma(b_2^* + w_2^* \times \hat{x}_1))\right)$

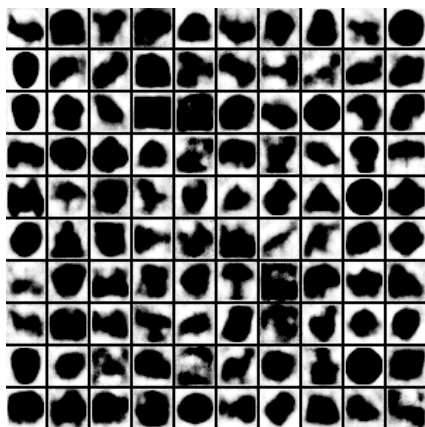
The generated sample is: $\hat{\mathbf{x}} = [\hat{x}_1, \hat{x}_2]^T$

Subsection 4

Application



(a) Dataset samples



(b) Generated samples

Figure: FVSBN performance over Caltech 101 dataset [7]

Section 5

Latent Variable Models with Likelihood Training

Evidence Lower Bound (ELBO)

ELBO

ELBO is a tractable lower bound to data log-likelihood as:

$$\text{ELBO}(\mathbf{x}_k) \leq \log p_{\theta}(\mathbf{x}_k)$$

Thus the following problem can be replaced for generative modeling:

$$\theta^* = \operatorname{argmax}_{\theta} \frac{1}{K} \sum_{k=1}^K \text{ELBO}(\mathbf{x}_k)$$

Models Using ELBO

- Variational AutoEncoder
- Diffusion probabilistic models

Change-of-Variable Technique

Consider the following two assumptions:

- $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I}), \mathbf{z} \in \mathbb{R}^D$
- $\mathbf{x} = \mathbf{f}_\theta(\mathbf{z})$ where $\mathbf{f}(\cdot) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is a nonlinear invertible function implemented by a DNN

Then $p_\theta(\mathbf{x}_i)$ can be calculated analytically using change-of-variable technique.

Models Using Change-of-Variable Technique

- Normalizing Flow

Section 6

Summary

Model	Density	Sampling	Training	Latent	Architecture
Autoregressive	Exact	Slow	MLE*	-	RNN/Transformer
GAN	-	Fast	JSD**	\mathbb{R}^L	Generator/Discriminator
VAE	Lower Bound	Fast	ELBO	\mathbb{R}^L	Encoder/Decoder
Normalizing Flow	Exact	Slow/Fast	MLE	\mathbb{R}^D	Invertible
Diffusion	Lower Bound	Slow	ELBO	\mathbb{R}^D	Encoder/Decoder

*Maximum Likelihood Estimation

**Jenson-Shanon Divergence

Table: Summary of Deep Generative Models

References I



Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu,
“Wavenet: A generative model for raw audio,”
arXiv preprint arXiv:1609.03499, 2016.



Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.,
“Photorealistic text-to-image diffusion models with deep language understanding,”
Advances in Neural Information Processing Systems, vol. 35, pp. 36479–36494, 2022.



Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi,
“Palette: Image-to-image diffusion models,”
in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–10.



Stefano Ermon and Yang Song,
“Cs236: Deep generative models,” Fall 2021.



Alec Radford, Luke Metz, and Soumith Chintala,
“Unsupervised representation learning with deep convolutional generative adversarial networks,”
arXiv preprint arXiv:1511.06434, 2015.



Guangyuan Zhang, Xiaoping Rui, Stefan Poslad, Xianfeng Song, Yonglei Fan, and Bang Wu, “A method for the estimation of finely-grained temporal spatial human population density distributions based on cell phone call detail records,” *Remote Sensing*, vol. 12, no. 16, pp. 2572, 2020.



Zhe Gan, Ricardo Henao, David Carlson, and Lawrence Carin, “Learning deep sigmoid belief networks with data augmentation,” in *Artificial Intelligence and Statistics*. PMLR, 2015, pp. 268–276.