

# Introduction to Deep Generative Modeling

## COMPSCI 589 - Summer 2024

Sajjad Amini

Manning College of Information and Computer Sciences  
University of Massachusetts Amherst

# Disclaimer

- The financial examples and data presented are for illustrative purposes only.

- All images in this presentation were generated using ChatGPT unless otherwise cited.
- Each image has been created to visually enhance the topics discussed and provide illustrative support.
- For images not generated by ChatGPT, sources are cited directly in the title.

# Contents

- 1 Intuition
- 2 Concept
- 3 Approaches
  - Autoregressive Modeling
  - Variational Autoencoder
  - Generative Adversarial Nets
  - Diffusion Models
- 4 Extention to Conditional Generation
- 5 Applications
- 6 Deep Autoregressive Models

# Section 1

## Intuition

# Investment Challenge



Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge



Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

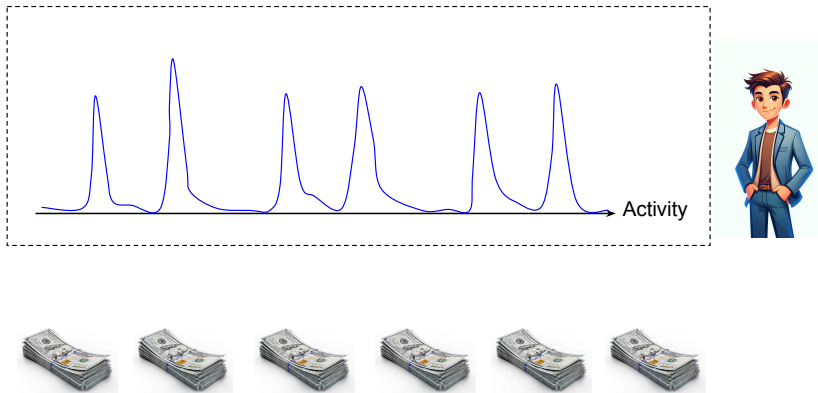


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

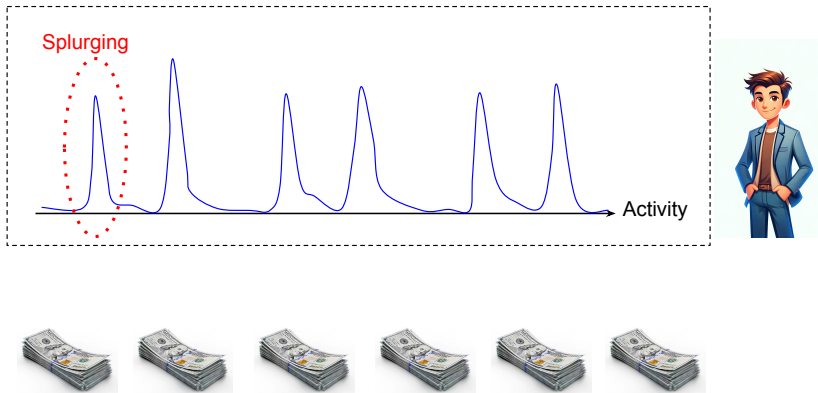


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

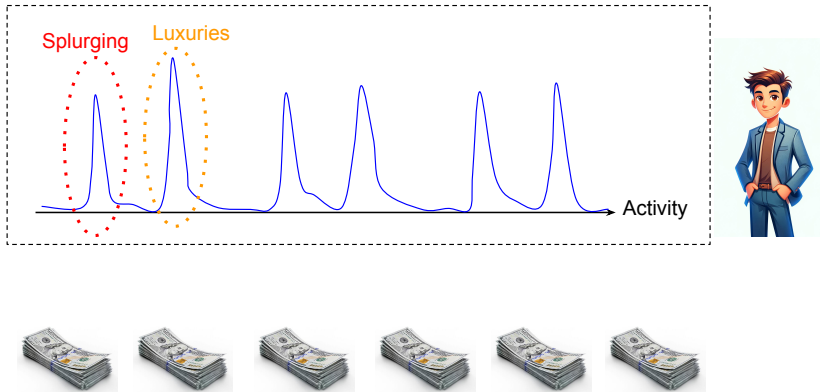


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

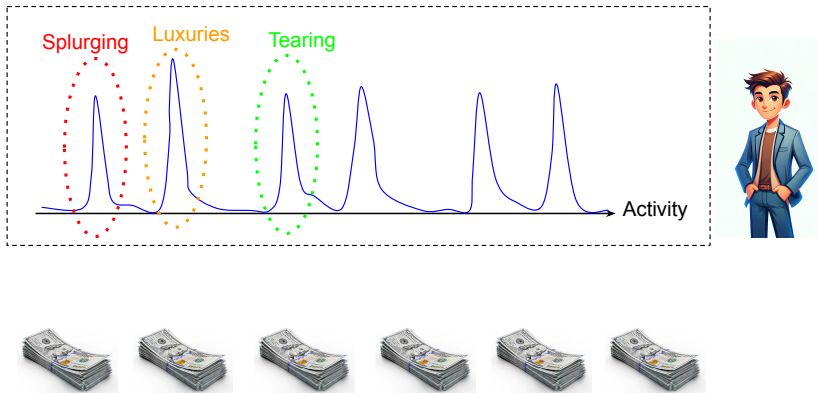


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

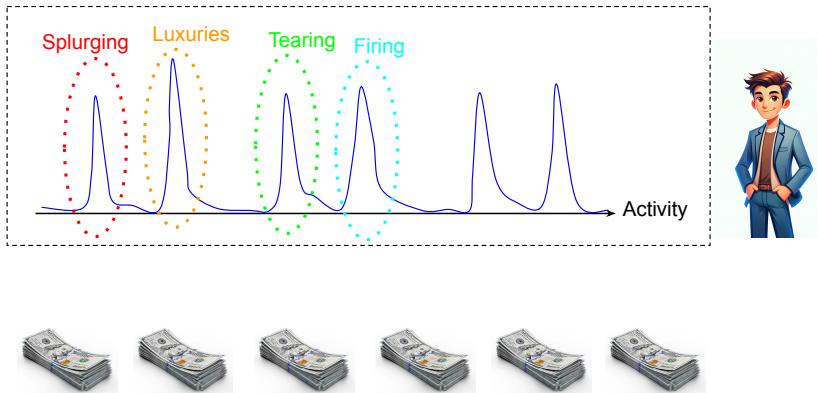


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

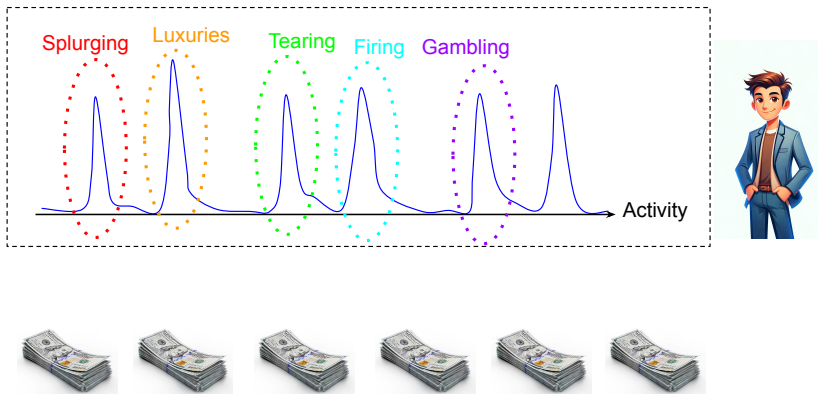


Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge

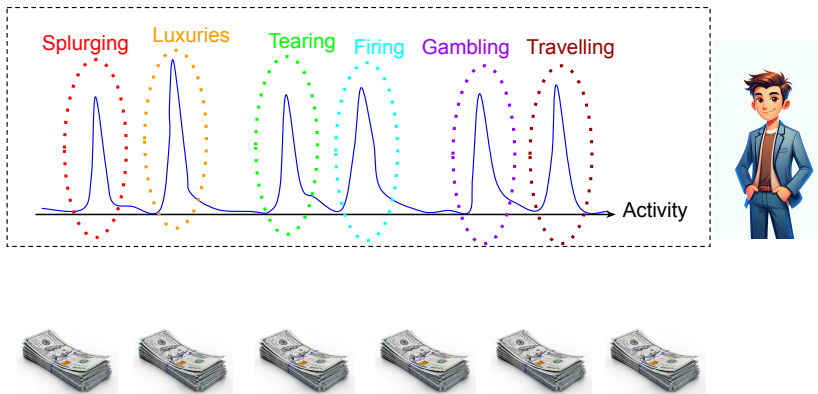
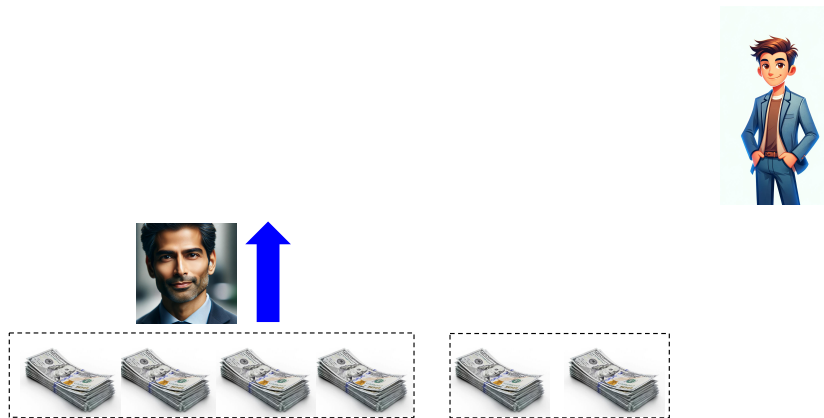


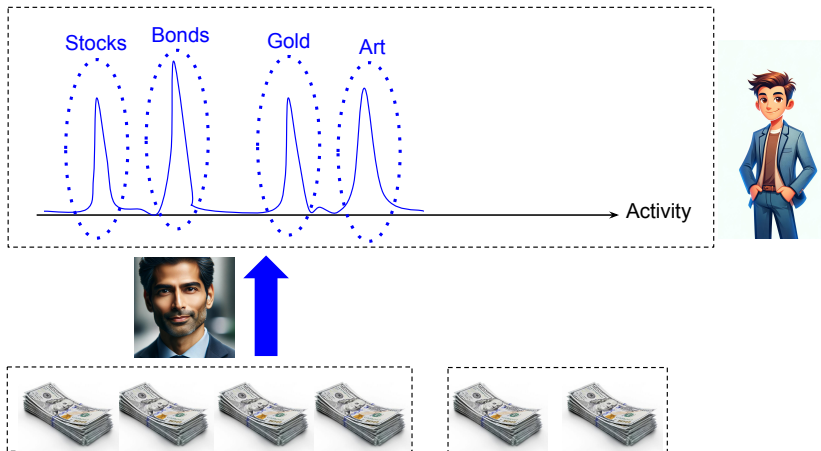
Figure: Investment Challenge (Budget: \$1M, Divesting is not allowed)

# Investment Challenge



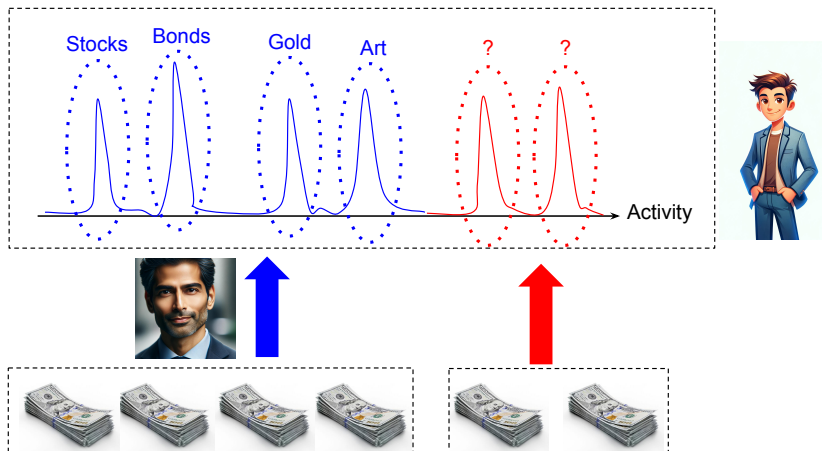
**Figure:** Investment challenge with help of investor (Budget: \$1M, Divesting is not allowed)

# Investment Challenge



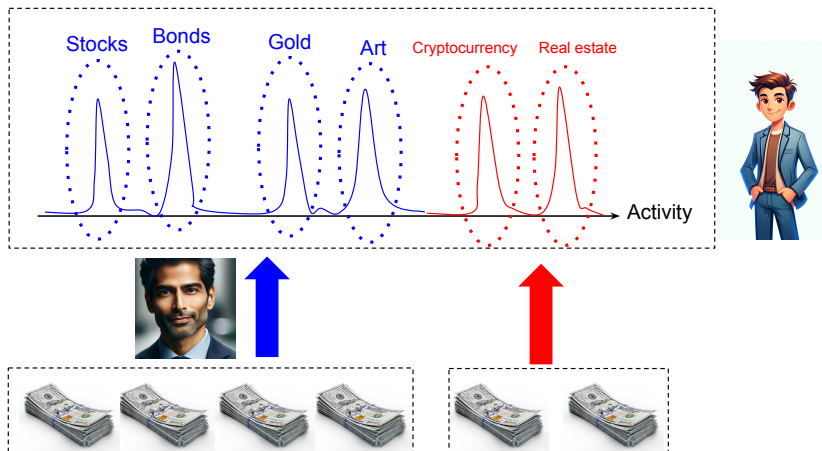
**Figure:** Investment challenge with help of investor (Budget: \$1M, Divesting is not allowed)

# Investment Challenge



**Figure:** Investment challenge with help of investor (Budget: \$1M, Divesting is not allowed)

# Investment Challenge



**Figure:** Investment challenge with help of investor (Budget:  $\$1M$ , Divesting is not allowed)

# Investment Challenge

Axioms of Investment Challenge	
Limited Budget	

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


Axioms of Investment Challenge	
Limited Budget	

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


Axioms of Investment Challenge	
Limited Budget	
Positivity	

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


Axioms of Investment Challenge	
Limited Budget	
Positivity	INVEST  DIVEST

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


Axioms of Investment Challenge	
Limited Budget	
Positivity	INVEST <del>DEVEST</del>

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


	Axioms of Investment Challenge	Axioms of Probability
Limited Budget		
Positivity	INVEST <del>DEVEST</del>	

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


	Axioms of Investment Challenge	Axioms of Probability
Limited Budget		$\int_x p(x) dx = 1$
Positivity	INVEST <del>DEVEST</del>	

Figure: From Axioms of our challenge to Axioms of probability

# Investment Challenge


	Axioms of Investment Challenge	Axioms of Probability
Limited Budget		$\int_{\mathcal{X}} p(\mathbf{x}) d\mathbf{x} = 1$
Positivity	INVEST <del>DISINVEST</del>	$p(\mathbf{x}) \geq 0$

Figure: From Axioms of our challenge to Axioms of probability

## Section 2

### Concept

# Parametric Probability Density Function (PDF)

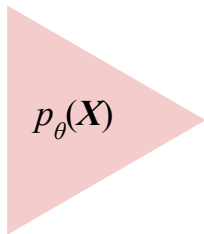


Figure: Your new budget is your parametric PDF

# Parametric Probability Density Function (PDF)

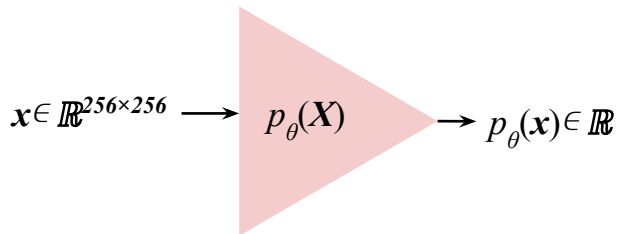


Figure: Your new budget is your parametric PDF

# Parametric Probability Density Function (PDF)

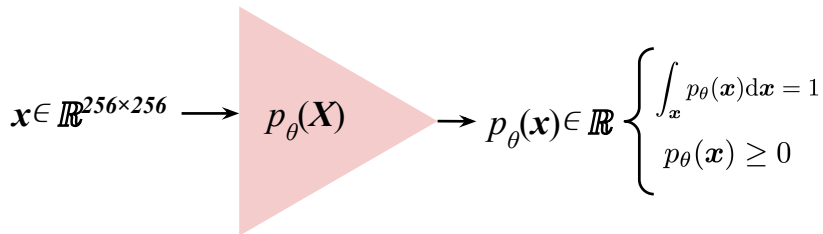


Figure: Your new budget is your parametric PDF

# Parametric Probability Density Function (PDF)

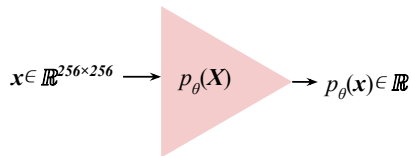


Figure: Your new budget is your parametric PDF

# Parametric Probability Density Function (PDF)

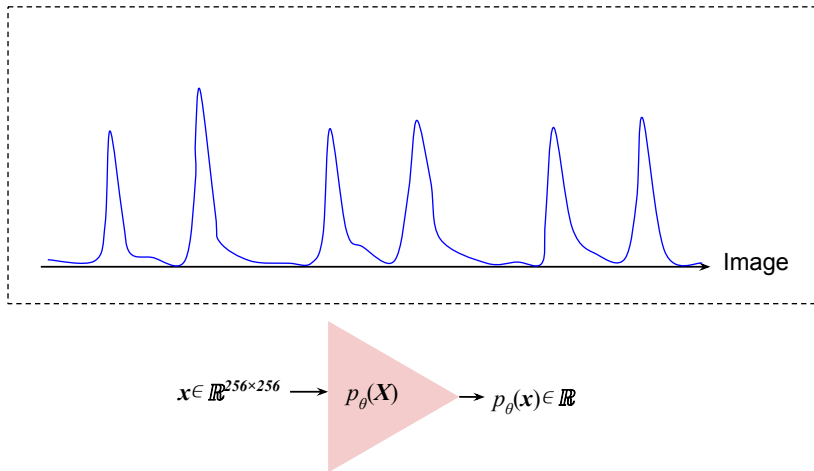


Figure: Your new budget is your parametric PDF

# Parametric Probability Density Function (PDF)

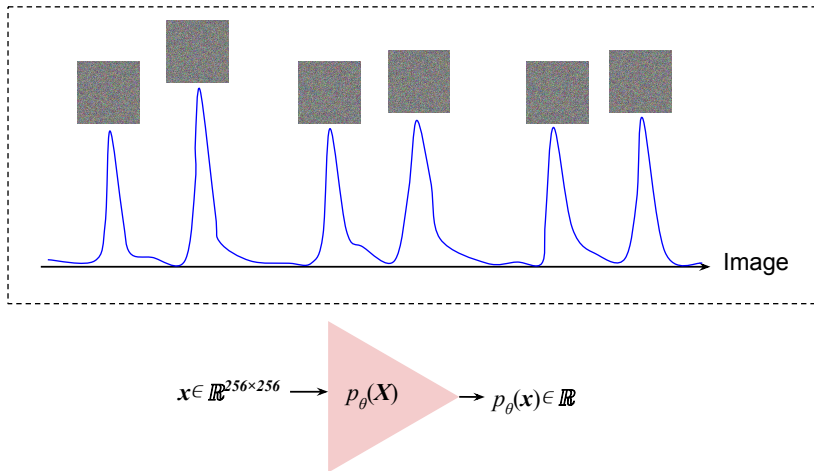


Figure: Your new budget is your parametric PDF

# Learning Rooms!

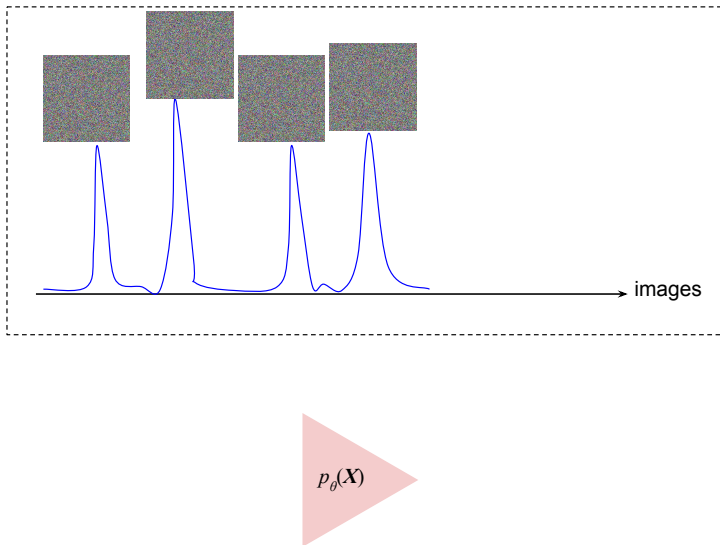


Figure: Learning to represent rooms

# Learning Rooms!

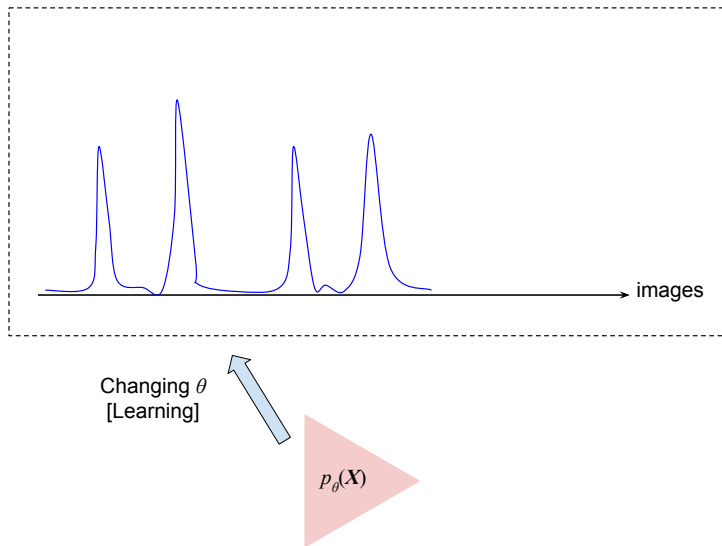


Figure: Learning to represent rooms

# Learning Rooms!

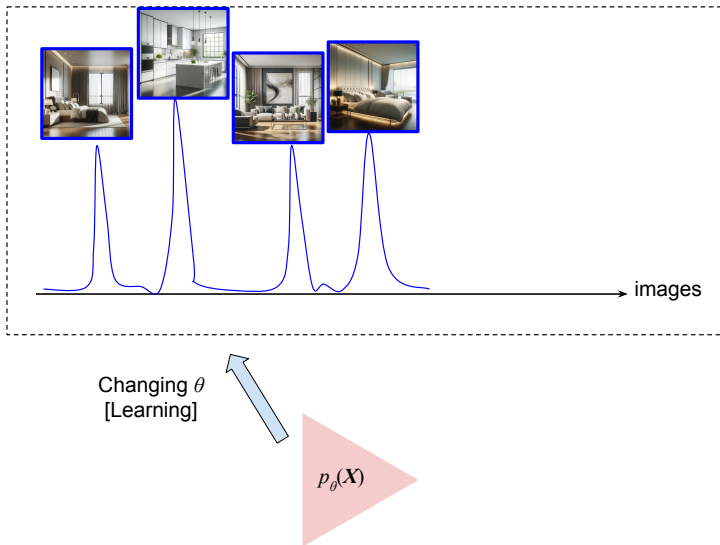


Figure: Learning to represent rooms

# Learning Rooms!

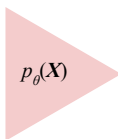
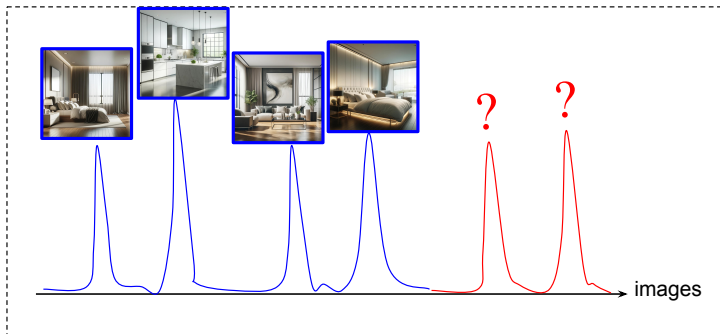


Figure: Learning to represent rooms

# Learning Rooms!

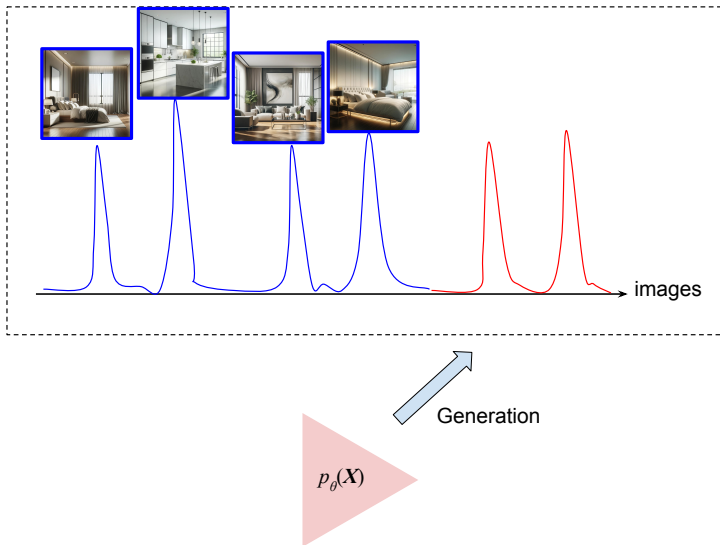


Figure: Learning to represent rooms

# Learning Rooms!

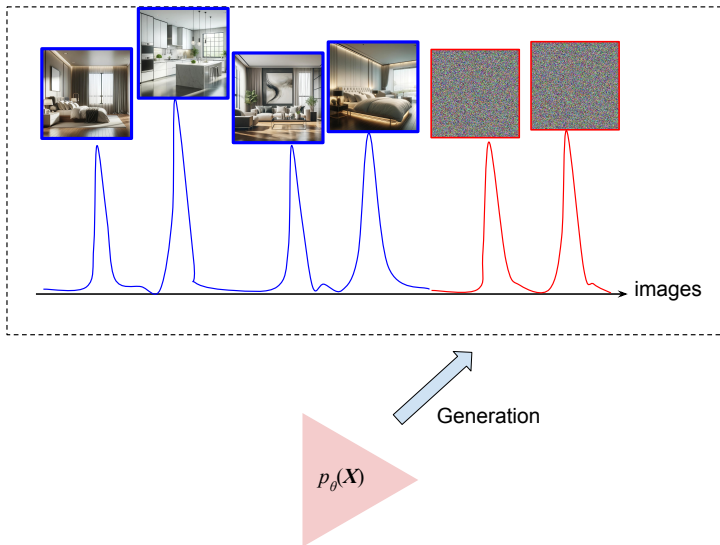


Figure: Learning to represent rooms

# Learning Rooms!

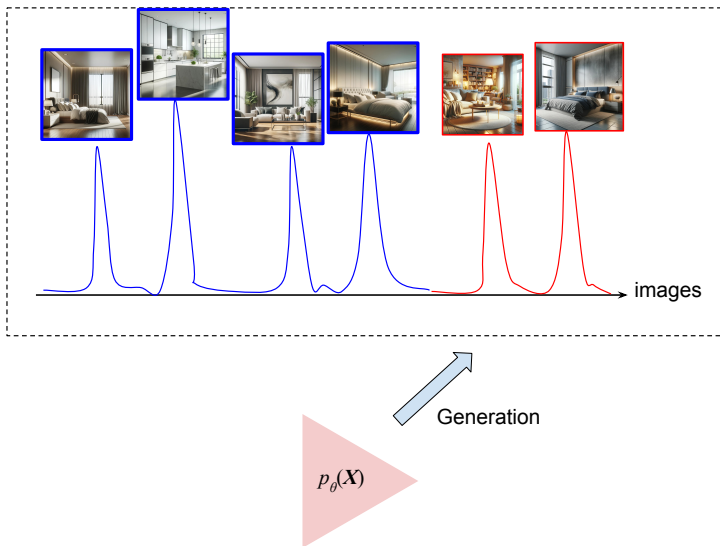


Figure: Learning to represent bedrooms

## Section 3

# Approaches

## Subsection 1

### Autoregressive Modeling

# Autoregressive Modeling

*"You can generate data if you can predict its future given its past!"*

# Language Modeling Using Autoregressive Models

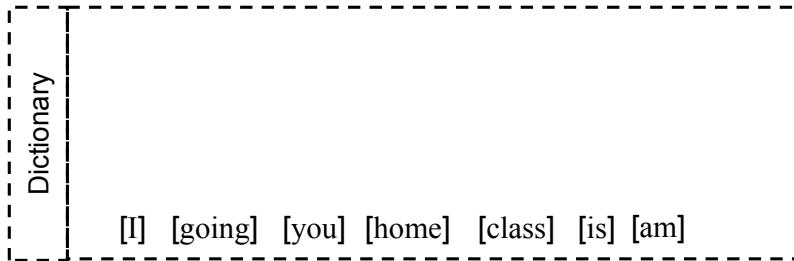


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

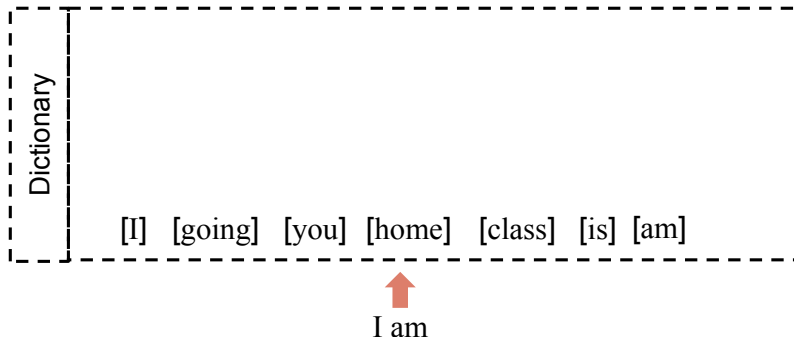


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

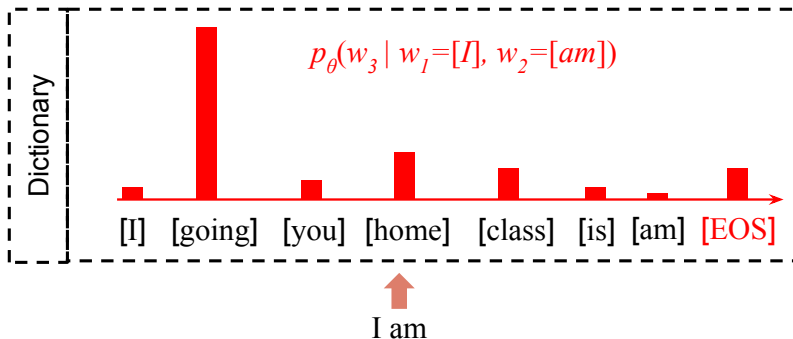


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

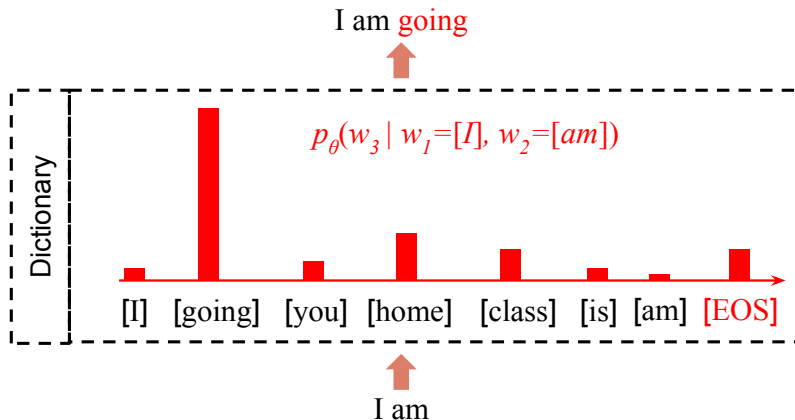


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

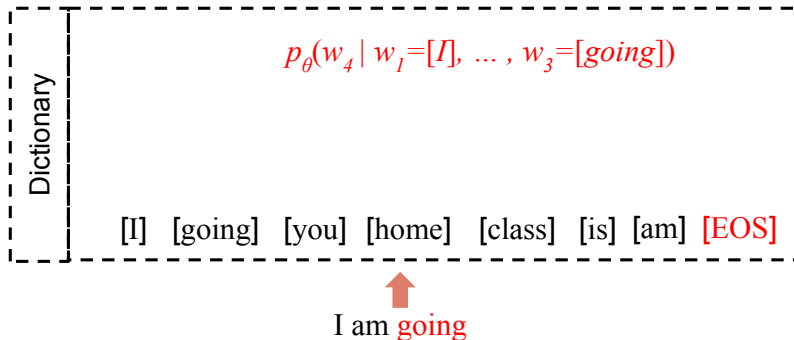


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

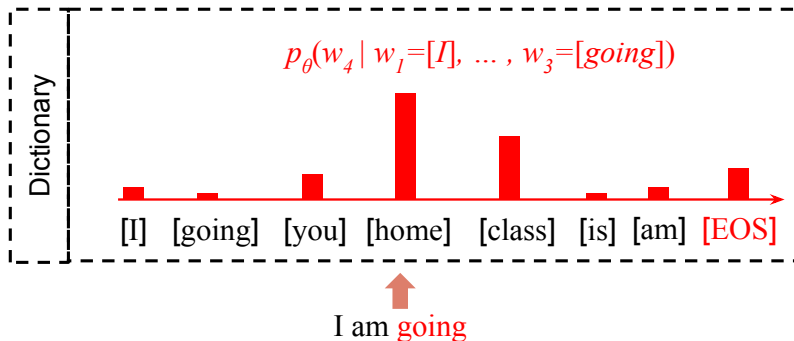


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

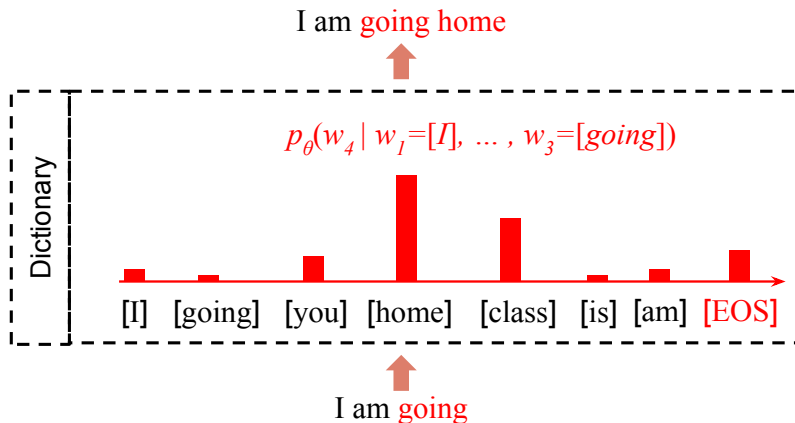


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

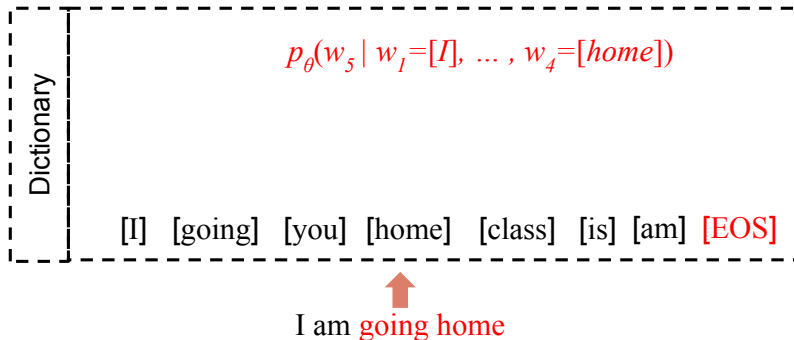


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

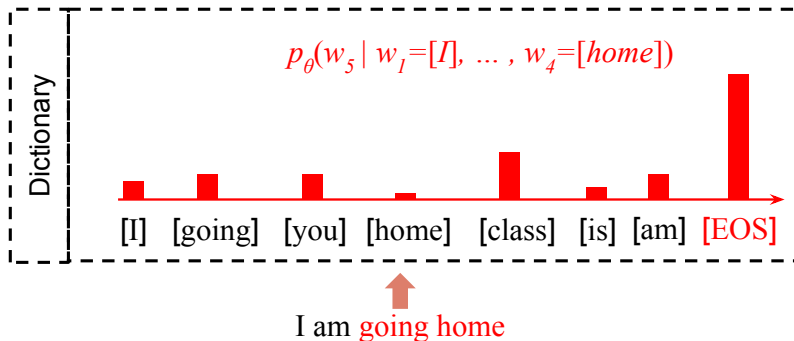


Figure: Generating the remaining part of a sentence

# Language Modeling Using Autoregressive Models

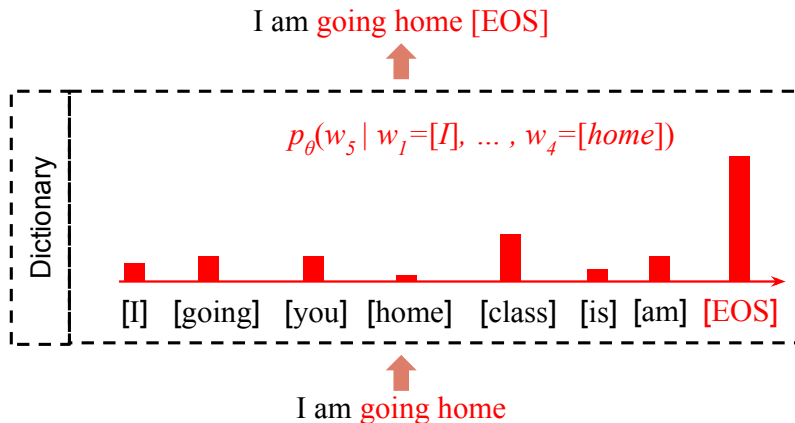


Figure: Generating the remaining part of a sentence

# Scaling to ChatGPT

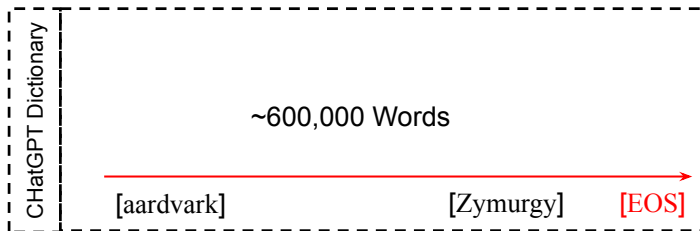
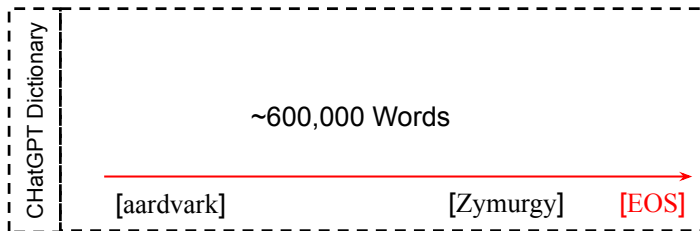


Figure: ChatGPT built on top of an Autoregressive model

# Scaling to ChatGPT

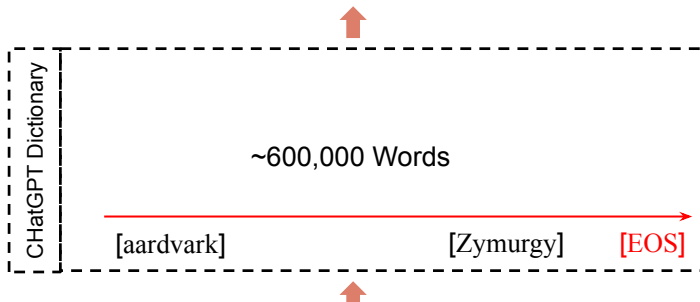


George has three brothers and one sister. How many people are in his family, including his mother and father?

Figure: ChatGPT built on top of an Autoregressive model

# Scaling to ChatGPT

George has three brothers and one sister. How many people are in his family, including his mother and father? George has three brothers and one sister, making a total of five children. Including his mother and father, there are seven people in George's family.



George has three brothers and one sister. How many people are in his family, including his mother and father?

Figure: ChatGPT built on top of an Autoregressive model

## Subsection 2

### Variational Autoencoder

*"You can generate data if you can compress it efficiently!"*

# Variational Autoencoders

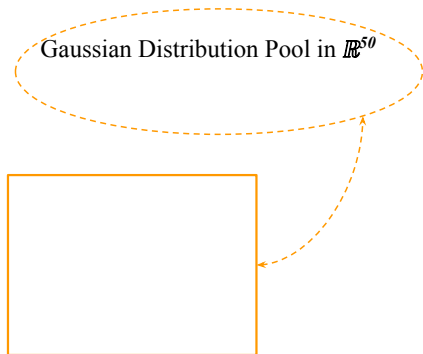


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

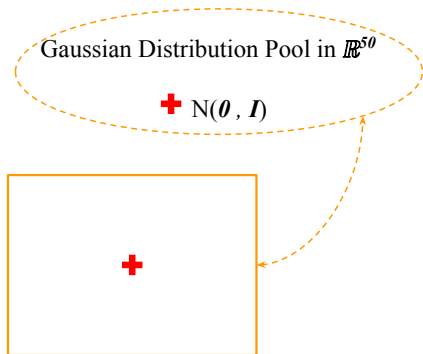


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders



$$x \in \mathbb{R}^{256 \times 256}$$

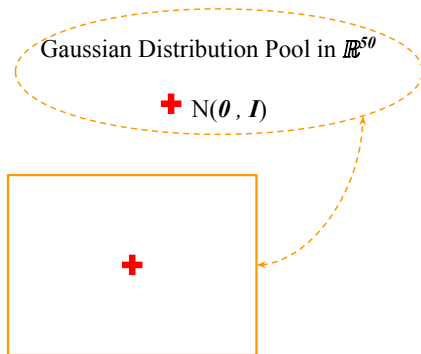


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

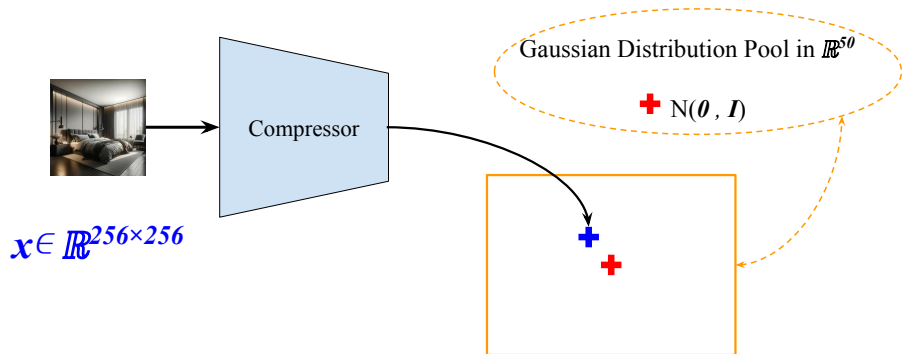


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

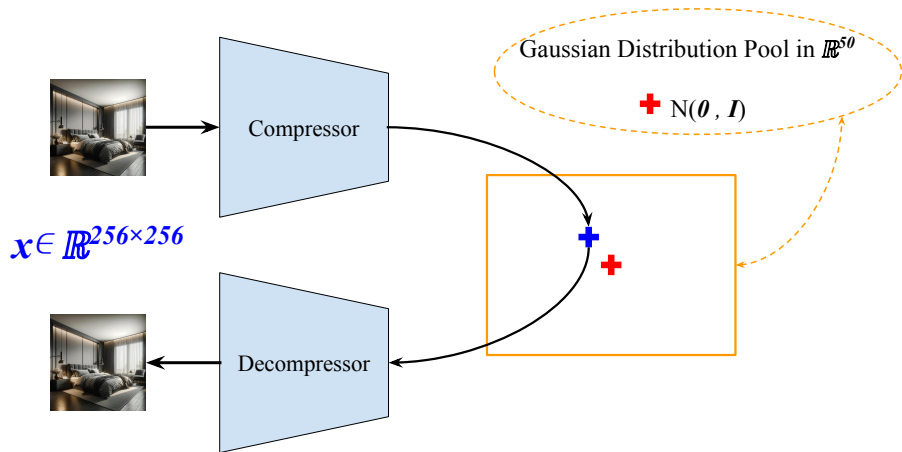


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

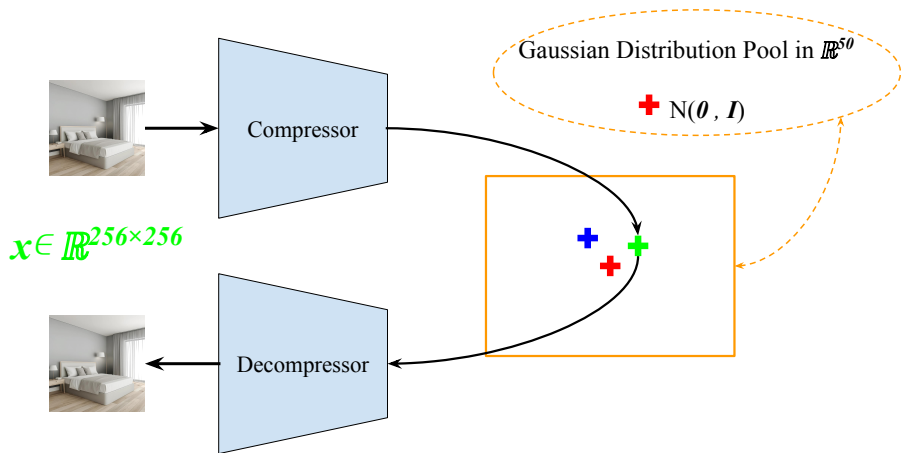


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

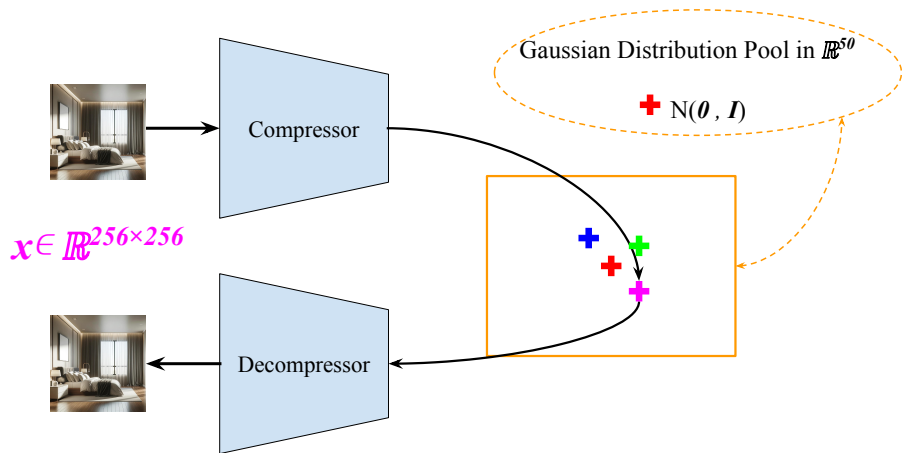


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

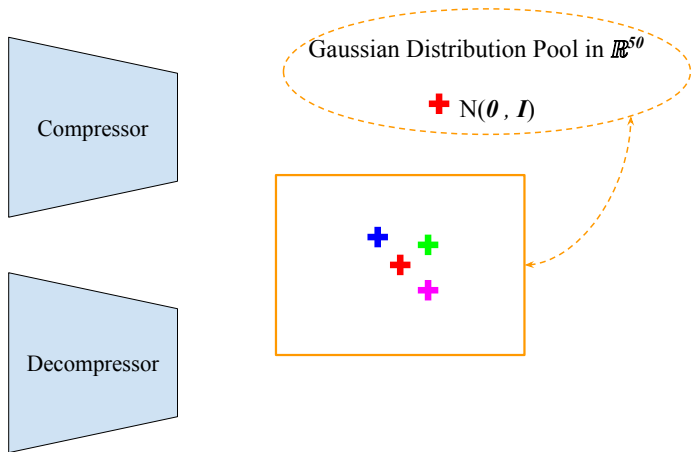


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

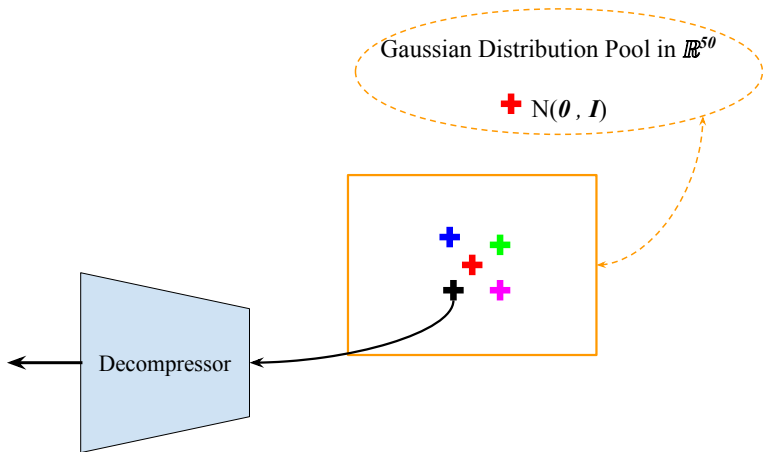


Figure: Compression learning as a method of generative modeling

# Variational Autoencoders

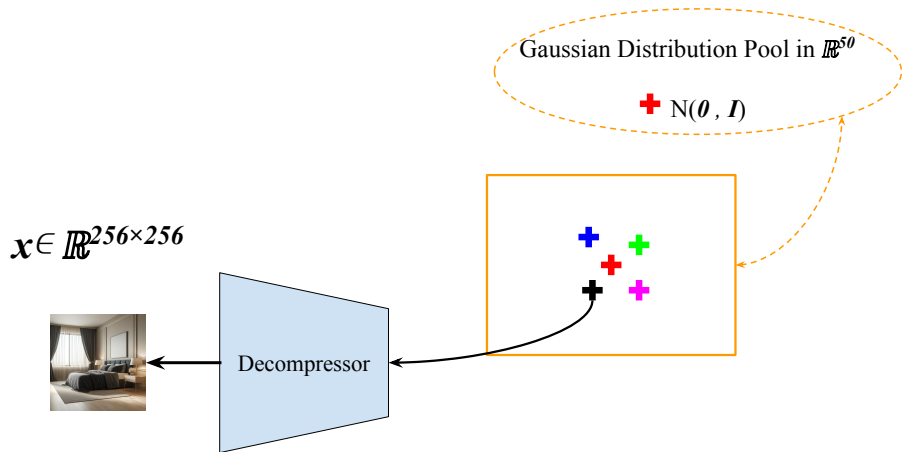


Figure: Compression learning as a method of generative modeling

## Subsection 3

### Generative Adversarial Nets

*"Good generated samples are those that are indistinguishable from the real ones!"*

# Generative Adversarial Nets

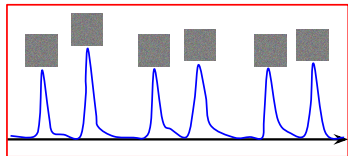


Figure: Using an Inspector [Discriminator] to detect generation

# Generative Adversarial Nets

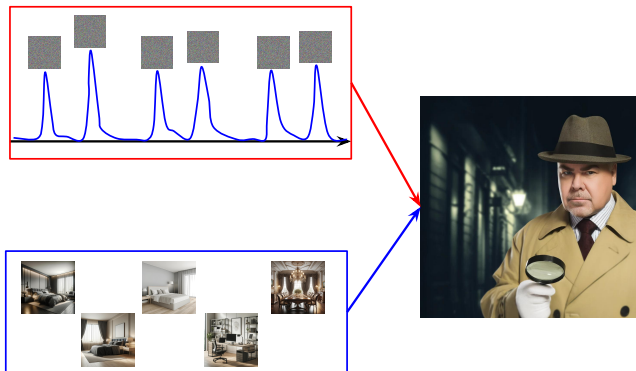


Figure: Using an Inspector [Discriminator] to detect generation

# Generative Adversarial Nets

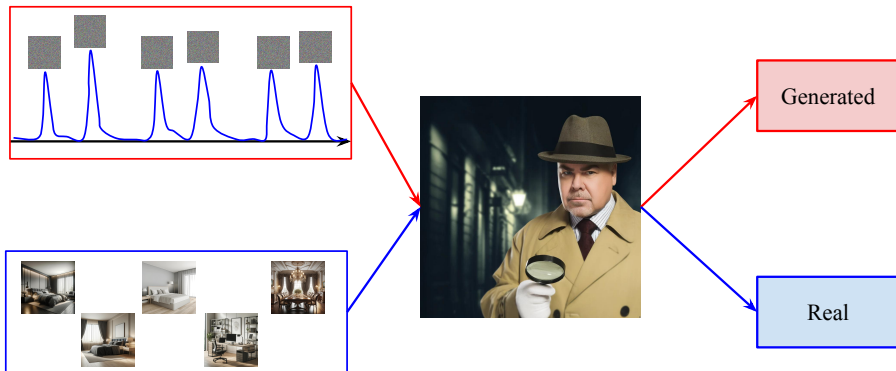


Figure: Using an Inspector [Discriminator] to detect generation

# Generative Adversarial Nets

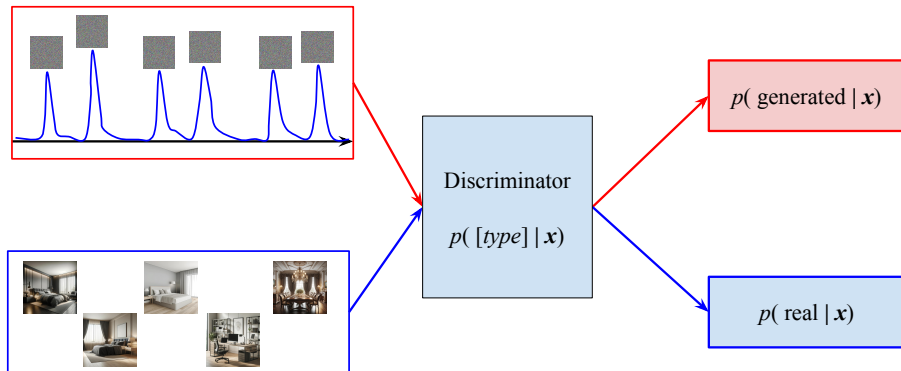


Figure: Using an Inspector [Discriminator] to detect generation

# Generative Adversarial Nets

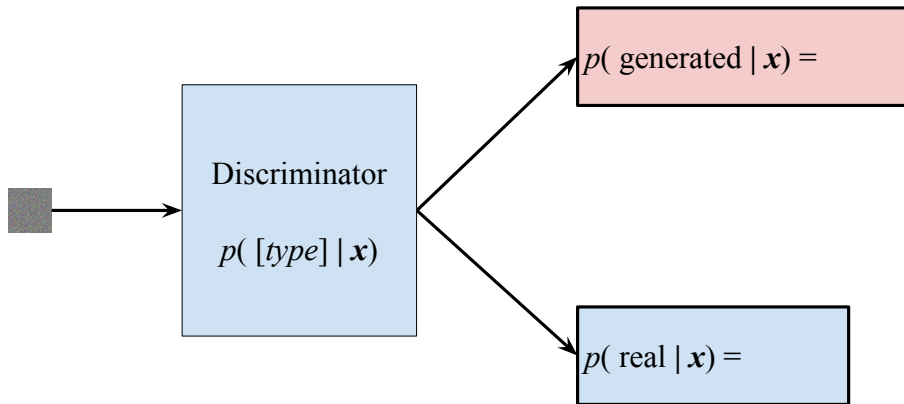


Figure: Examining the Discriminator

# Generative Adversarial Nets

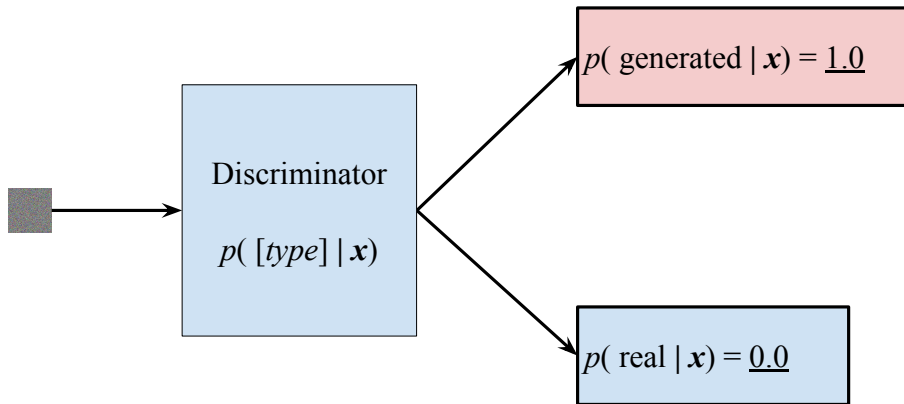


Figure: Examining the Discriminator

# Generative Adversarial Nets

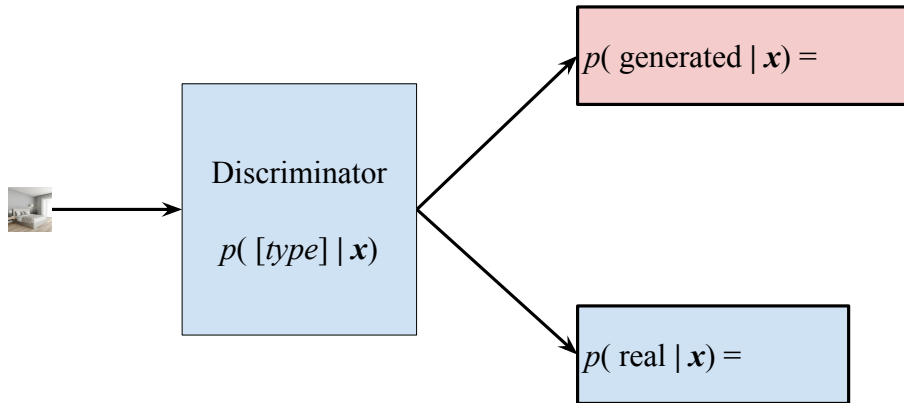


Figure: Examining the Discriminator

# Generative Adversarial Nets

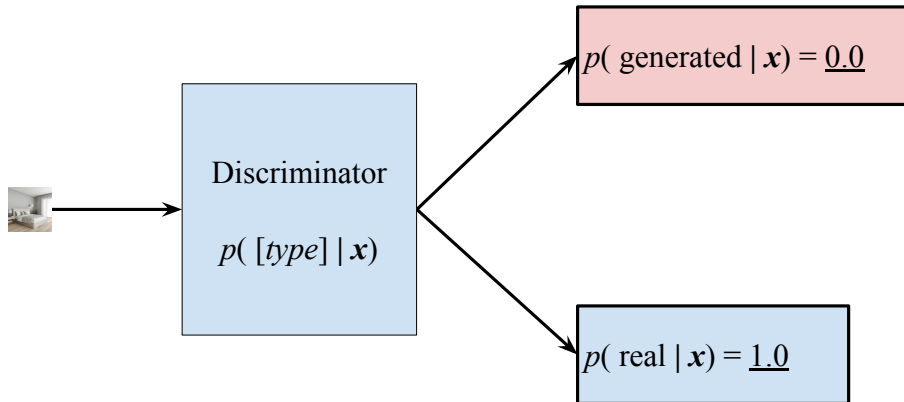


Figure: Examining the Discriminator

# Generative Adversarial Nets

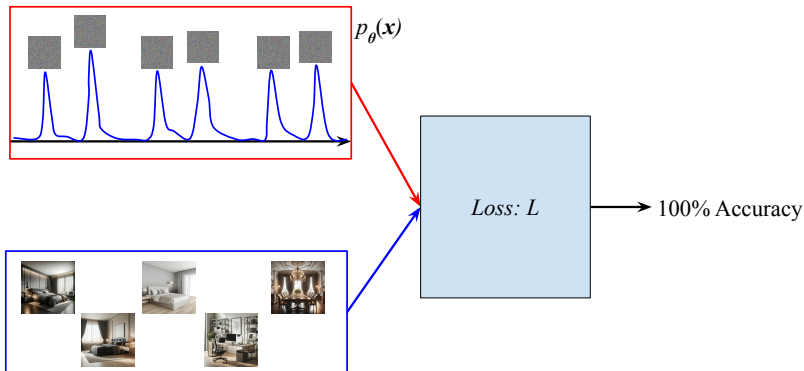


Figure: Updating generation

# Generative Adversarial Nets

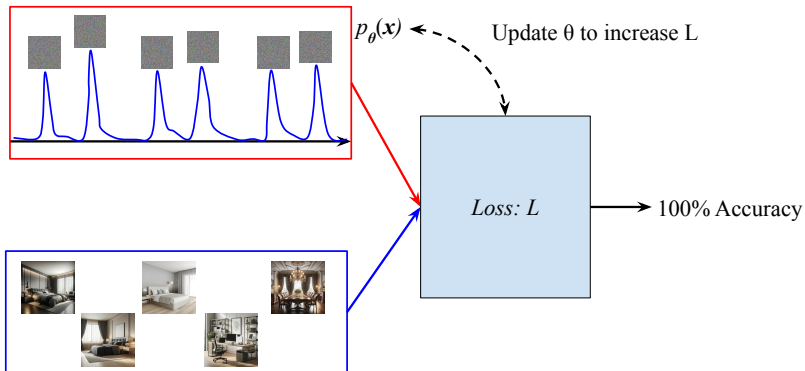


Figure: Updating generation

# Generative Adversarial Nets

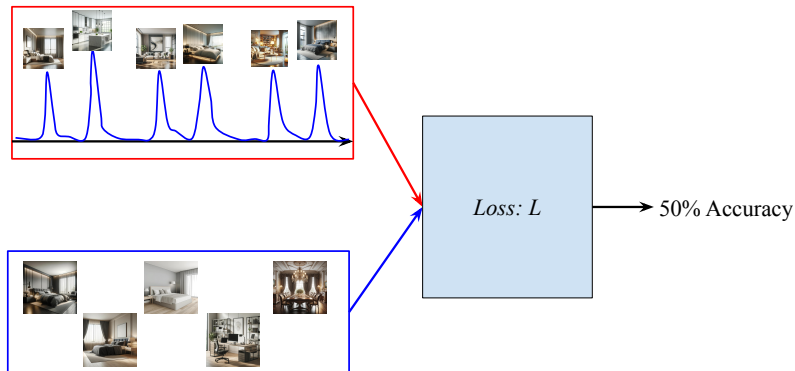


Figure: Updating generation

## Subsection 4

### Diffusion Models

*"You can generate data if you can denoise it"*



$\sigma$

Figure: Denoiser module

# Diffusion Models Denoiser

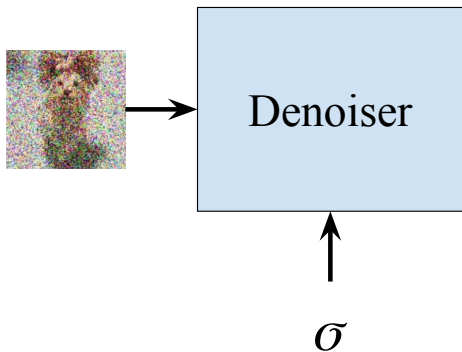


Figure: Denoiser module

# Diffusion Models Denoiser

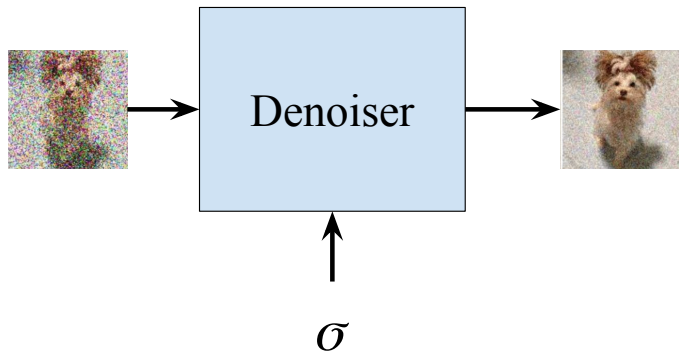


Figure: Denoiser module

# Diffusion Models Generation

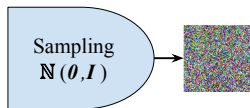


Figure: Generation using diffusion model (images source: [1])

# Diffusion Models Generation

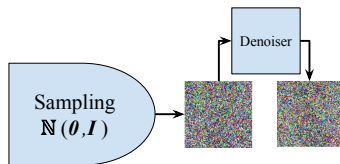


Figure: Generation using diffusion model (images source: [1])

# Diffusion Models Generation

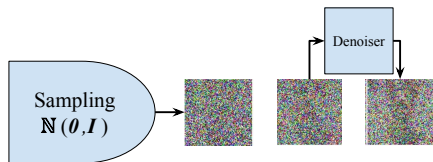


Figure: Generation using diffusion model (images source: [1])

# Diffusion Models Generation

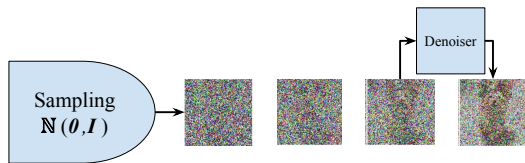


Figure: Generation using diffusion modeld (images source: [1])

# Diffusion Models Generation

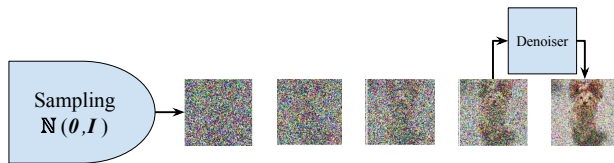


Figure: Generation using diffusion modeld (images source: [1])

# Diffusion Models Generation

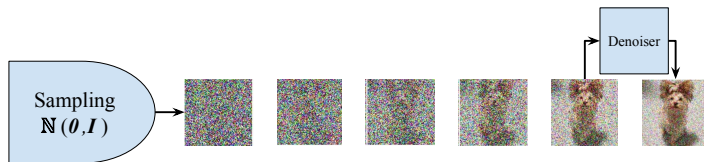


Figure: Generation using diffusion modeld (images source: [1])

# Diffusion Models Generation

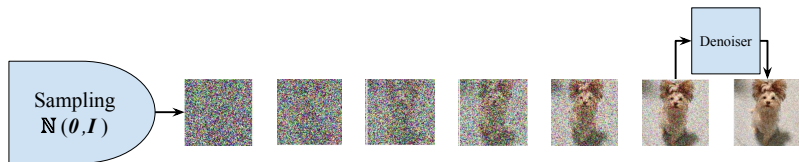


Figure: Generation using diffusion modeld (images source: [1])

# Diffusion Models Generation

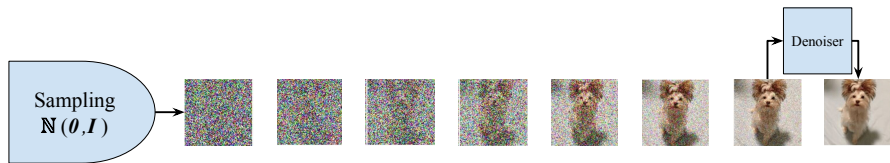


Figure: Generation using diffusion modeld (images source: [1])

## Section 4

### Extention to Conditional Generation

# Learning Conditional Distributions

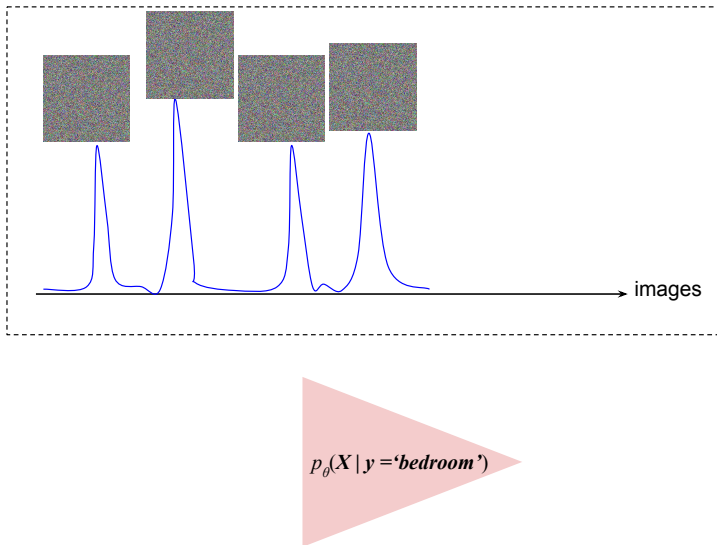


Figure: Learning to represent bedrooms

# Learning Conditional Distributions

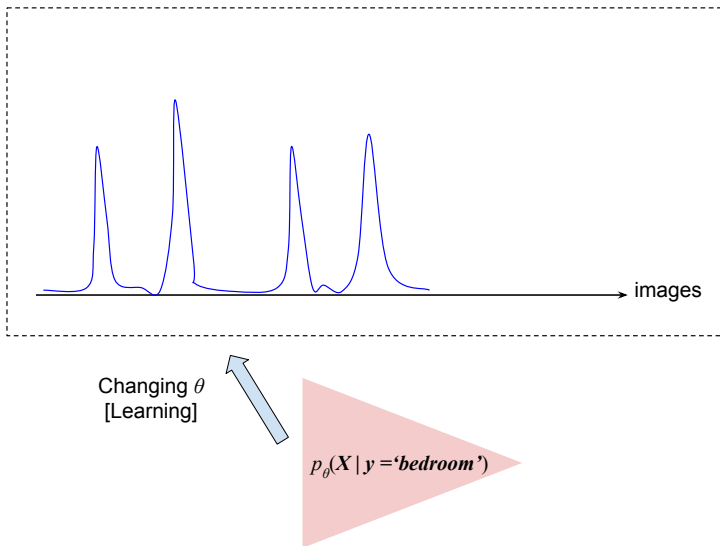


Figure: Learning to represent bedrooms

# Learning Conditional Distributions

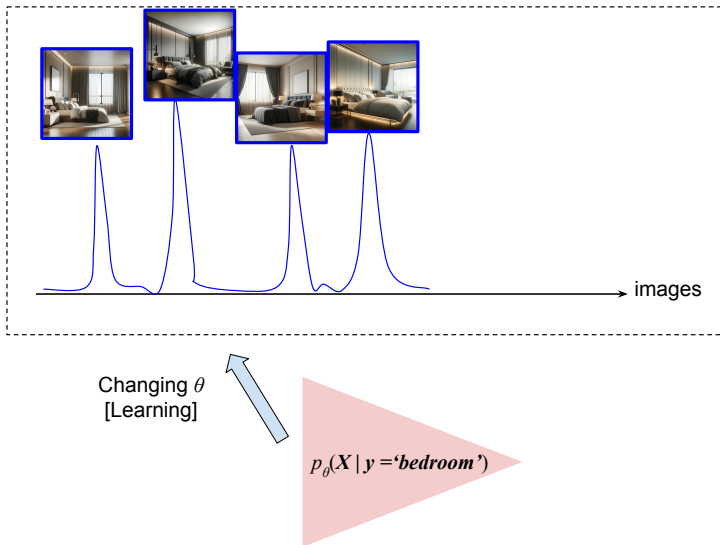
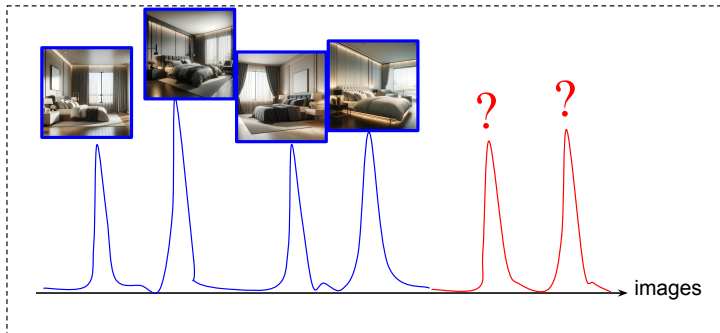


Figure: Learning to represent bedrooms

# Learning Conditional Distributions



$$p_{\theta}(X|y=\text{'bedroom'})$$

Figure: Learning to represent bedrooms

# Learning Conditional Distributions

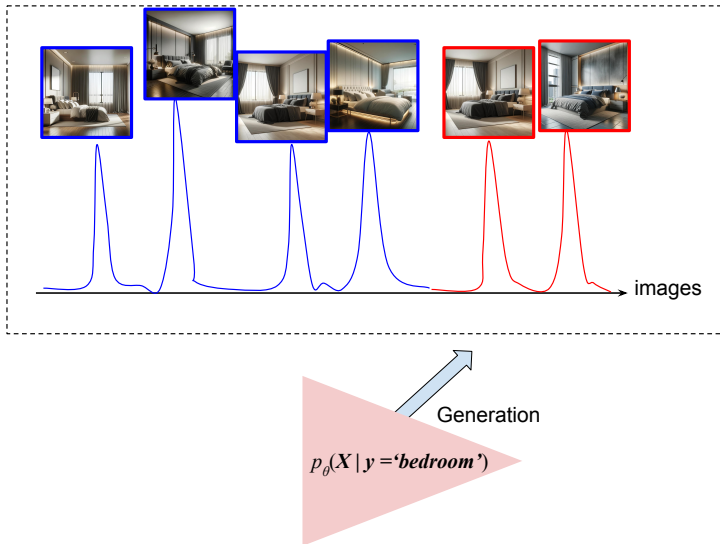


Figure: Learning to represent bedrooms

# Learning Conditional Distributions

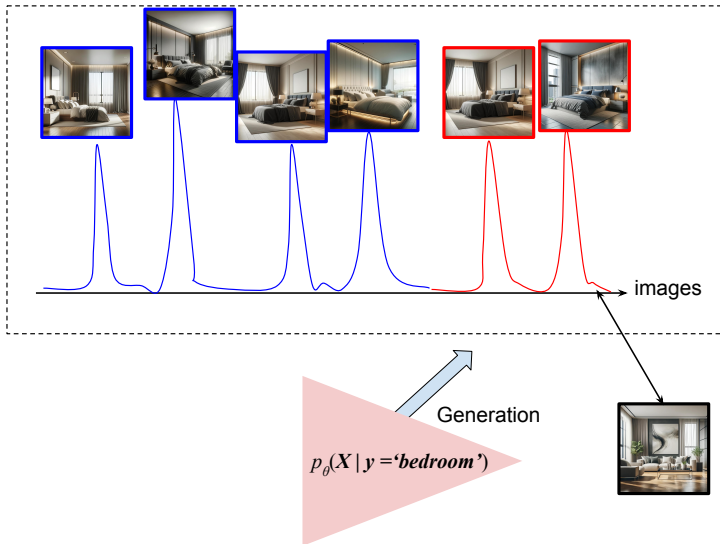


Figure: Learning to represent bedrooms

## Section 5

# Applications

# Text-to-Speech Models

## Text-to-Speech Models

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

# Text-to-Speech Models

## Text-to-Speech Models

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

## Real-World Sample

Listen to the following speech synthesis (source: [2])

$$\mathbf{y} = \begin{array}{c} \text{"A single Wavenet can} \\ \text{capture the characteristics of many} \\ \text{different speakers with equal fidelity,} \\ \text{not it's fast."} \end{array} \xrightarrow{\text{Sampling } p(\mathbf{x}|\mathbf{y})} \mathbf{x} = \boxed{\text{Play}}$$

# Text-to-Image Models

## Text-to-Image Models

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

# Text-to-Image Models

## Text-to-Image Models

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



Figure:  $\mathbf{x}$  for  $\mathbf{y}$  = “Teddy bears swimming at the Olympics 400m Butterfly event.”  
(source: [?])

# Image-to-Image Translation

## Image Colorization

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

# Image-to-Image Translation

## Image Colorization

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



(a)  $\mathbf{y}$



(b)  $\mathbf{x}$



(c) Ground truth

Figure: Image colorization (source: [3])

## Image Inpainting

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

# Image-to-Image Translation

## Image Inpainting

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



(a)  $\mathbf{y}$



(b)  $\mathbf{x}$



(c) Ground truth

Figure: Image inpainting (source: [3])

## Image Uncropping

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

# Image-to-Image Translation

## Image Uncropping

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



(a)  $\mathbf{y}$



(b)  $\mathbf{x}$



(c) Ground truth

Figure: Image uncropping (source: [3])

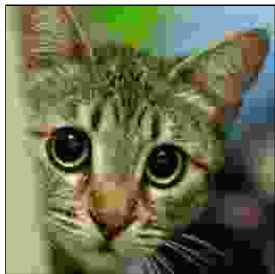
## Image Restoration

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

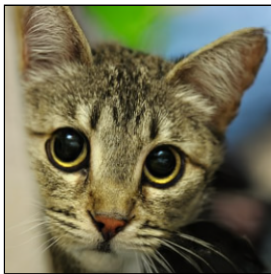
# Image-to-Image Translation

## Image Restoration

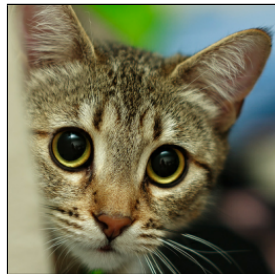
$$p(\mathbf{x}|\mathbf{y}) : \begin{cases} \mathbf{x} : \text{A clean image} \\ \mathbf{y} : \text{A degraded image} \end{cases}$$



(a)  $\mathbf{y}$



(b)  $\mathbf{x}$



(c) Ground truth

Figure: Image restoration (source: [3])

## Section 6

# Deep Autoregressive Models

# Logistic Regression Model

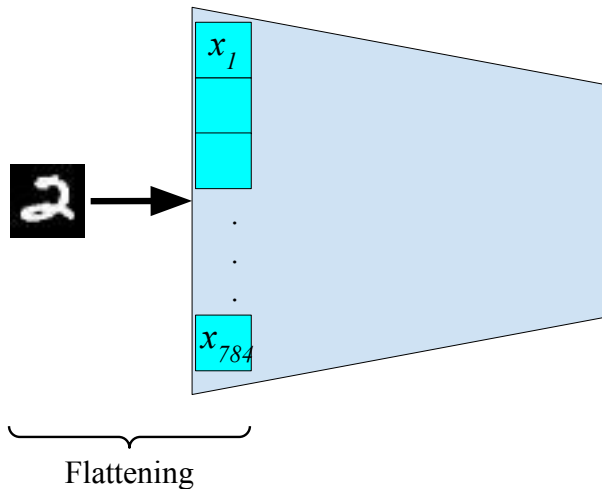


Figure: Logistic regression steps

# Logistic Regression Model

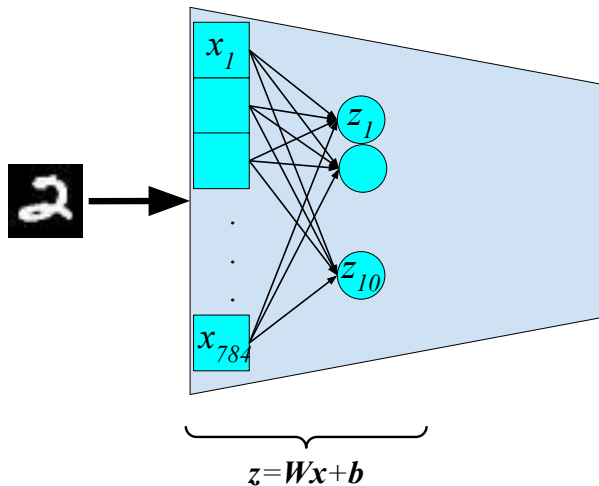


Figure: Logistic regression steps

# Logistic Regression Model

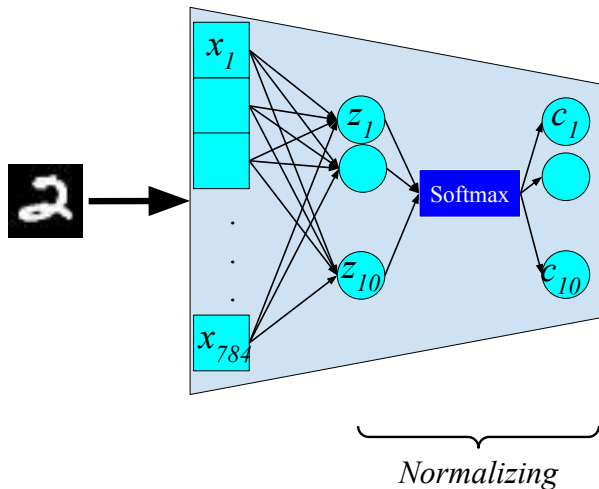
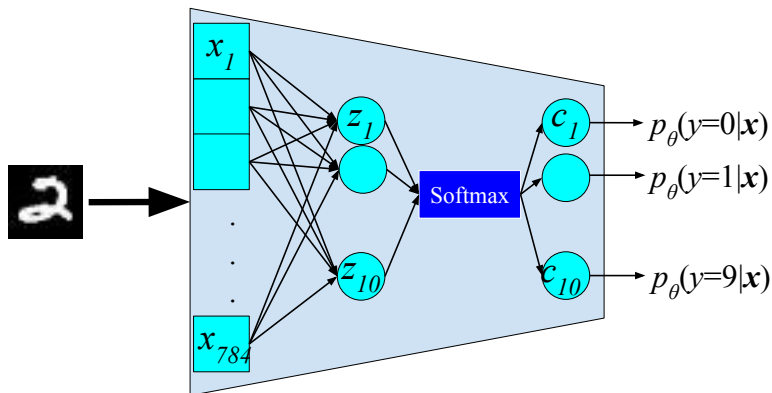


Figure: Logistic regression steps

# Logistic Regression Model



$$\underbrace{\qquad\qquad\qquad}_{p_{\theta}(y|\mathbf{x}) = \text{Cat}(y;\mathbf{C})}$$

Figure: Logistic regression steps

# Logistic Regression Model

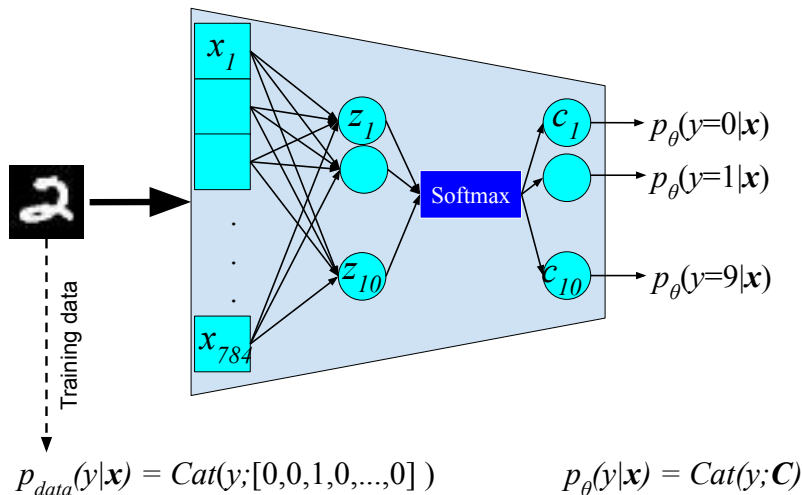


Figure: Logistic regression steps

# Logistic Regression Model

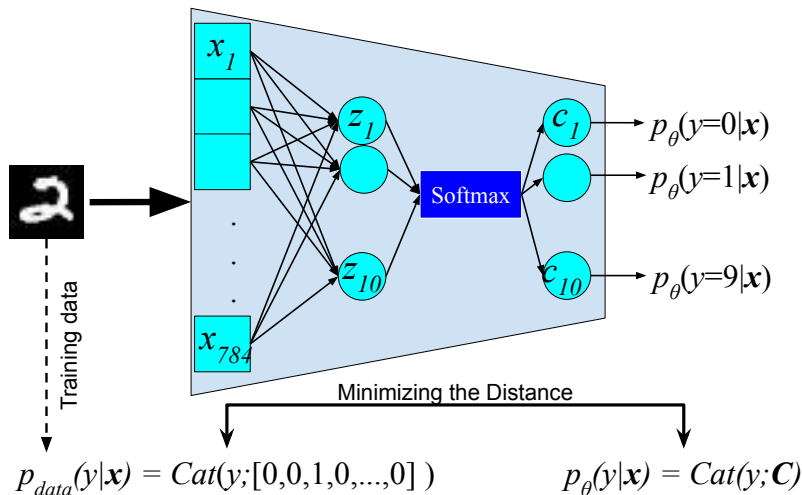


Figure: Logistic regression steps

## Distance Metric

One option for distance metric is:

## Distance Metric

One option for distance metric is:

$$L(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \text{KL} \left( p_{\text{data}}(y|\boldsymbol{x}) \parallel p_{\boldsymbol{\theta}}(y|\boldsymbol{x}) \right) \right]$$

## Distance Metric

One option for distance metric is:

$$\begin{aligned} L(\boldsymbol{\theta}) &= \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \text{KL} \left( p_{\text{data}}(y|\boldsymbol{x}) \parallel p_{\boldsymbol{\theta}}(y|\boldsymbol{x}) \right) \right] \\ &= \sum_{\boldsymbol{x}} p_{\text{data}}(\boldsymbol{x}) \left[ \sum_y p_{\text{data}}(y|\boldsymbol{x}) \log \frac{p_{\text{data}}(y|\boldsymbol{x})}{p_{\boldsymbol{\theta}}(y|\boldsymbol{x})} \right] \end{aligned}$$

## Distance Metric

One option for distance metric is:

$$\begin{aligned} L(\theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \text{KL} \left( p_{\text{data}}(y|\mathbf{x}) \parallel p_{\theta}(y|\mathbf{x}) \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \left[ \sum_y p_{\text{data}}(y|\mathbf{x}) \log \frac{p_{\text{data}}(y|\mathbf{x})}{p_{\theta}(y|\mathbf{x})} \right] \\ &= \underbrace{\sum_y \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}, y) \log p_{\text{data}}(y|\mathbf{x})}_{\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\text{data}}(y|\mathbf{x})]} \end{aligned}$$

## Distance Metric

One option for distance metric is:

$$\begin{aligned} L(\theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \text{KL} \left( p_{\text{data}}(y|\mathbf{x}) \parallel p_{\theta}(y|\mathbf{x}) \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \left[ \sum_y p_{\text{data}}(y|\mathbf{x}) \log \frac{p_{\text{data}}(y|\mathbf{x})}{p_{\theta}(y|\mathbf{x})} \right] \\ &= \underbrace{\sum_y \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}, y) \log p_{\text{data}}(y|\mathbf{x})}_{\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\text{data}}(y|\mathbf{x})]} - \underbrace{\sum_y \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}, y) \log p_{\theta}(\mathbf{x}|y)}_{\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\theta}(\mathbf{x}|y)]} \end{aligned}$$

## Distance Metric

One option for distance metric is:

$$\begin{aligned} L(\theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \text{KL} \left( p_{\text{data}}(y|\mathbf{x}) \parallel p_{\theta}(y|\mathbf{x}) \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \left[ \sum_y p_{\text{data}}(y|\mathbf{x}) \log \frac{p_{\text{data}}(y|\mathbf{x})}{p_{\theta}(y|\mathbf{x})} \right] \\ &= \underbrace{\sum_y \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}, y) \log p_{\text{data}}(y|\mathbf{x})}_{\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\text{data}}(y|\mathbf{x})]} - \underbrace{\sum_y \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}, y) \log p_{\theta}(\mathbf{x}|y)}_{\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\theta}(\mathbf{x}|y)]} \end{aligned}$$

While the second term is a function of your model parameters, the first one is independent of the selected Autoregressive model and thus can be omitted in optimization.

## Distance Metric

So:

$$\operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} -\mathbb{E}_{(\mathbf{x},y) \sim p_{\text{data}}(\mathbb{X},Y)} [\log p_{\boldsymbol{\theta}}(y|\mathbf{x})]$$

## Monte Carlo Estimation

Consider the following expectation:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

## Distance Metric

So:

$$\operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} -\mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\boldsymbol{\theta}}(y|\mathbf{x})]$$

## Monte Carlo Estimation

Consider the following expectation:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

Now assume that instead of  $p(\mathbb{X})$ , we just have access to  $N$  independent samples of random variable  $\mathbb{X}$  as  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

## Distance Metric

So:

$$\operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} -\mathbb{E}_{(\mathbf{x},y) \sim p_{\text{data}}(\mathbb{X},Y)} [\log p_{\boldsymbol{\theta}}(y|\mathbf{x})]$$

## Monte Carlo Estimation

Consider the following expectation:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] = \int p(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}$$

Now assume that instead of  $p(\mathbb{X})$ , we just have access to  $N$  independent samples of random variable  $\mathbb{X}$  as  $\mathbf{x}_1, \dots, \mathbf{x}_N$ . Then expectation can be approximated as:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbb{X})} [f(\mathbf{x})] \simeq \frac{1}{N} \sum_n f(\mathbf{x}_n)$$

## Optimization

Using Monte-Carlo estimation, we have the following optimization problem:

$$\begin{aligned}\boldsymbol{\theta}^* &= \operatorname{argmax}_{\boldsymbol{\theta}} -\mathbb{E}_{(\boldsymbol{x}, y) \sim p_{\text{data}}(\mathbb{X}, Y)} [\log p_{\boldsymbol{\theta}}(y|\boldsymbol{x})] \\ &\simeq \operatorname{argmax}_{\boldsymbol{\theta}} -\frac{1}{N} \sum_{i=1}^N \log p_{\boldsymbol{\theta}}(y_i|\boldsymbol{x}_i)\end{aligned}$$

# Sampling

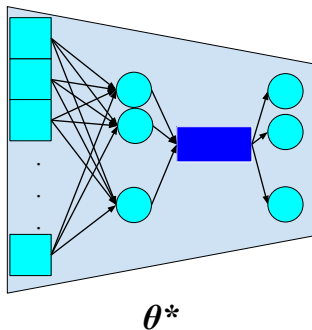


Figure: Sampling a trained model

# Sampling

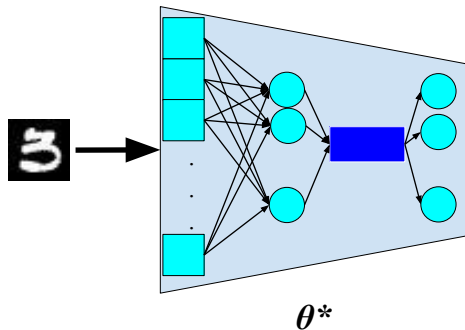


Figure: Sampling a trained model

# Sampling

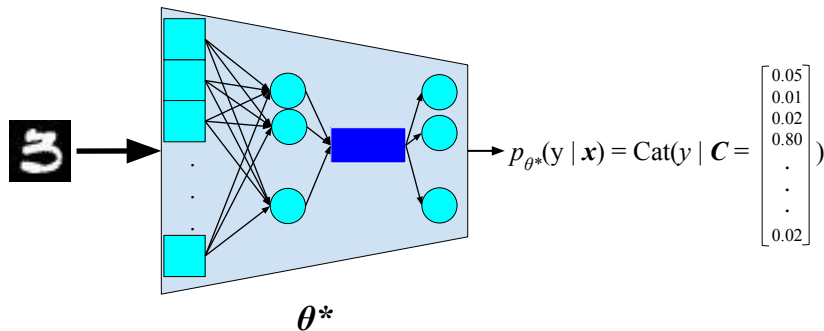


Figure: Sampling a trained model

# Sampling

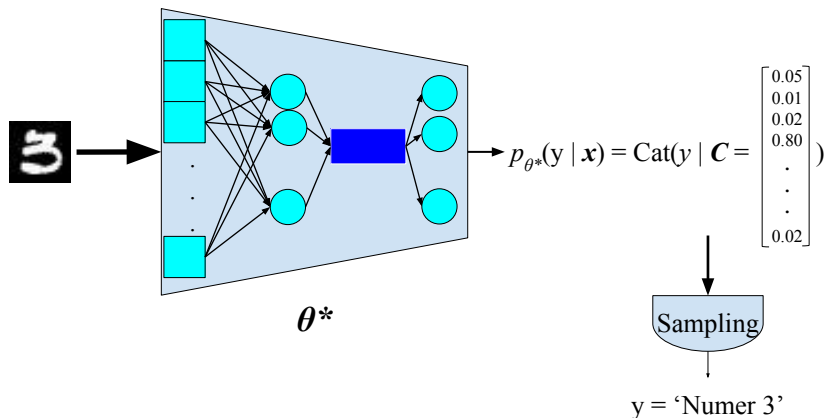


Figure: Sampling a trained model

# Modeling

## Generative Modeling

Assume we just have MNIST image  $\{\mathbf{x}_i\}_{i=1}^N$  without any label and we want to estimate generating distribution  $p(\mathbf{x})$  where  $\mathbf{x} \in \mathbb{R}^{784}$ .

## Challenge: High-dimensional Random Vector

In contrast to logistic regression where we model  $p_{\text{data}}(y|\mathbf{x})$  and  $y$  was a one-dimensional random variable, here  $\mathbf{x}$  is a high-dimensional random vector.

- ☞ It seems that we can't use logistic regression here.
- ☞ We can model each dimension separately because  $x_i \in \{0, 1, 2, \dots, 255\}$

## Chain Rule

Based on the chain rule, we have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_d|\mathbf{x}_{<d}) \dots p(x_D|\mathbf{x}_{<D}), \quad \mathbf{x}_{<d} \triangleq [x_1, \dots, x_{d-1}]^T$$

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times \boxed{p(x_i | \mathbf{x}_{<i})} \times \dots \times p(x_D | \mathbf{x}_{<D})$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times \boxed{p(x_i | \mathbf{x}_{<i})} \times \dots \times p(x_D | \mathbf{x}_{<D})$$

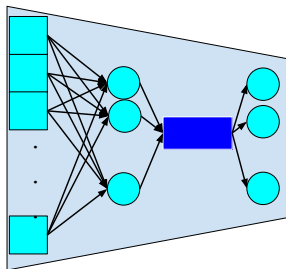
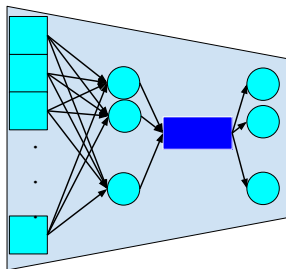


Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times \boxed{p(x_i | \mathbf{x}_{<i})} \times \dots \times p(x_D | \mathbf{x}_{<D})$$



$$\mathbf{W}_i, \mathbf{b}_i$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$

Figure: Using logistic regression for generative modeling

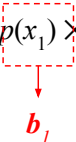
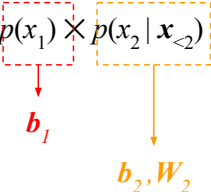
$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$


Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$


The diagram illustrates the decomposition of a joint probability distribution  $p(\mathbf{x})$  into a product of conditional distributions. The first term,  $p(x_1)$ , is highlighted with a red dashed box and a red arrow pointing to the parameter  $b_1$ . The second term,  $p(x_2 | \mathbf{x}_{<2})$ , is highlighted with an orange dashed box and an orange arrow pointing to the parameters  $b_2, W_2$ .

Figure: Using logistic regression for generative modeling

# Modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$

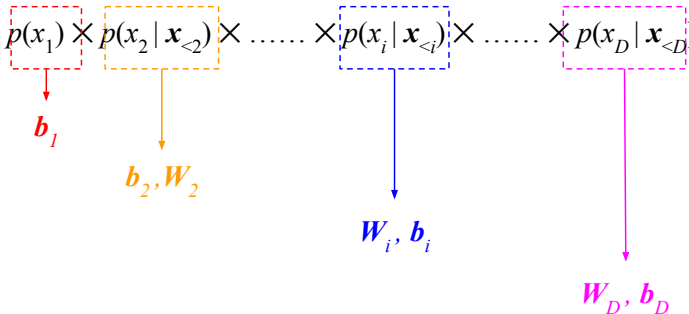
$\mathbf{b}_1$

$\mathbf{b}_2, \mathbf{W}_2$

$\mathbf{W}_i, \mathbf{b}_i$

Figure: Using logistic regression for generative modeling

# Modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$


The diagram illustrates the decomposition of a joint probability distribution  $p(\mathbf{x})$  into a product of conditional distributions. Each term in the product is enclosed in a dashed box of a different color. Arrows point from these boxes to the parameters of the corresponding logistic regression models:

- $p(x_1)$  (red dashed box) points to  $\mathbf{b}_1$  (red).
- $p(x_2 | \mathbf{x}_{<2})$  (orange dashed box) points to  $\mathbf{b}_2, \mathbf{W}_2$  (orange).
- $p(x_i | \mathbf{x}_{<i})$  (blue dashed box) points to  $\mathbf{W}_i, \mathbf{b}_i$  (blue).
- $p(x_D | \mathbf{x}_{<D})$  (magenta dashed box) points to  $\mathbf{W}_D, \mathbf{b}_D$  (magenta).

Figure: Using logistic regression for generative modeling

# Modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$

$\mathbf{b}_1$

$\mathbf{b}_2, \mathbf{W}_2$

$\mathbf{W}_i, \mathbf{b}_i$

$\mathbf{W}_D, \mathbf{b}_D$

$$x_i \in \{0, 1, \dots, 255\} \Rightarrow \begin{cases} \mathbf{b}_i \in \mathbb{R}^{256} \\ \mathbf{W}_i \in \mathbb{R}^{256 \times i} \end{cases} \quad \forall \quad 1 \leq i \leq D$$

Figure: Using logistic regression for generative modeling

# Modeling

$$p(\mathbf{x}) = p(x_1) \times p(x_2 | \mathbf{x}_{<2}) \times \dots \times p(x_i | \mathbf{x}_{<i}) \times \dots \times p(x_D | \mathbf{x}_{<D})$$

$\mathbf{b}_1$

$\mathbf{b}_2, \mathbf{W}_2$

$\mathbf{W}_i, \mathbf{b}_i$

$\mathbf{W}_D, \mathbf{b}_D$

$$\theta = \{ \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_i, \mathbf{b}_i, \dots, \mathbf{W}_D, \mathbf{b}_D \}$$

Figure: Using logistic regression for generative modeling

# Distance Metric

## Distance Metric

We want to compare two distributions  $p_{\text{data}}$  and  $p_{\theta}$ , thus we can use KL divergence as:

$$L(\theta) = \text{KL}(p_{\text{data}} \| p_{\theta}) =$$

# Distance Metric

## Distance Metric

We want to compare two distributions  $p_{\text{data}}$  and  $p_{\theta}$ , thus we can use KL divergence as:

$$L(\theta) = \text{KL}(p_{\text{data}} \| p_{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right]$$

## Distance Metric

We want to compare two distributions  $p_{\text{data}}$  and  $p_{\theta}$ , thus we can use KL divergence as:

$$\begin{aligned} L(\theta) &= \text{KL}(p_{\text{data}} \| p_{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \end{aligned}$$

## Distance Metric

We want to compare two distributions  $p_{\text{data}}$  and  $p_{\theta}$ , thus we can use KL divergence as:

$$\begin{aligned} L(\theta) &= \text{KL}(p_{\text{data}} \| p_{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \end{aligned}$$

Using above definition, we know  $L(\theta) = 0$  iff  $p_{\theta}(\mathbb{X}) = p_{\text{data}}(\mathbb{X})$ . We can rewrite  $L(\theta)$  as:

$$L(\theta) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})]$$

## Distance Metric

We want to compare two distributions  $p_{\text{data}}$  and  $p_{\theta}$ , thus we can use KL divergence as:

$$\begin{aligned} L(\theta) &= \text{KL}(p_{\text{data}} \| p_{\theta}) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ \log \left( \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right) \right] \\ &= \sum_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \end{aligned}$$

Using above definition, we know  $L(\theta) = 0$  iff  $p_{\theta}(\mathbb{X}) = p_{\text{data}}(\mathbb{X})$ . We can rewrite  $L(\theta)$  as:

$$L(\theta) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})] - \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})]$$

Because the first term on the right-hand side is independent of  $\theta$ , we have:

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

# From KL divergence to Model Likelihood

## Model Likelihood

We see:

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

Thus:

## Model Likelihood

We see:

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

Thus:

- Desirable situation is when  $p_{\theta}(\mathbb{X})$  assign high probability to probable regions in  $p_{\text{data}}(\mathbb{X})$

## Model Likelihood

We see:

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

Thus:

- Desirable situation is when  $p_{\theta}(\mathbb{X})$  assign high probability to probable regions in  $p_{\text{data}}(\mathbb{X})$
- We have yet a problem: No access to  $p_{\text{data}}$

# From KL divergence to Model Likelihood

## Model Likelihood

We see:

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

Thus:

- Desirable situation is when  $p_{\theta}(\mathbb{X})$  assign high probability to probable regions in  $p_{\text{data}}(\mathbb{X})$
- We have yet a problem: No access to  $p_{\text{data}}$
- $\mathbb{H}(p_{\text{data}}(\mathbb{X})) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} [\log p_{\text{data}}(\mathbf{x})]$  is the maximum accessible objective value where  $\mathbb{H}(p_{\text{data}}(\mathbb{X}))$  is the *entropy* defined as:

$$\mathbb{H}(p_{\text{data}}(\mathbb{X})) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\text{data}}(\mathbf{x})]$$

# Model Likelihood Estimation

## Model Likelihood Estimation

We are interested in solving the following problem:

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

but we don't have access to  $p_{\text{data}}$  and instead, we have access to independent samples from the distribution  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ .

# Model Likelihood Estimation

## Model Likelihood Estimation

We are interested in solving the following problem:

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

but we don't have access to  $p_{\text{data}}$  and instead, we have access to independent samples from the distribution  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ .

## Solution via Monte Carlo Estimate

Using the Monte Carlo estimate we have:

$$\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} [\log p_{\theta}(\mathbf{x})] \simeq \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(\mathbf{x}_n)$$

Thus:

$$\theta^* = \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(\mathbf{x}_n)$$

# Thank You!

Thank you for your attention!

*Do you have any questions or comments?*

## **Contact Information**

Sajjad Amini

Email: [samini@umass.edu](mailto:samini@umass.edu)

# References I



Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole,  
“Score-based generative modeling through stochastic differential equations,”  
*arXiv preprint arXiv:2011.13456*, 2020.



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, 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.