

Lecture 3: Deep Autoregressive Models

Deep Generative Models

Sajjad Amini

Department of Electrical Engineering
Sharif University of Technology

Contents

- 1 Fully Visible Sigmoid Belief Network
- 2 Neural Autoregressive Density Estimation
- 3 Masked Autoencoder for Distribution Estimation
- 4 Recurrent Neural Networks
- 5 Learning
- 6 Examples
- 7 Conditional Generation
- 8 Applications
- 9 Summary

Section 1

Fully Visible Sigmoid Belief Network

Binary Logistic Regression

Problem Statement

The original Binary MNIST dataset is $\{(\mathbf{x}_n, \bar{y}_n)\}_{n=1}^N$ where:

$$\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ \bar{y} \in \{0, 1, \dots, 9\} \end{cases}$$

Binary Logistic Regression

Problem Statement

The original Binary MNIST dataset is $\{(\mathbf{x}_n, \bar{y}_n)\}_{n=1}^N$ where:

$$\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ \bar{y} \in \{0, 1, \dots, 9\} \end{cases}$$

Assume we define a new problem where the labels are binarized as:

$$y_n = \begin{cases} 0 & \text{if } \bar{y}_n < 5 \\ 1 & \text{Otherwise} \end{cases}$$

Binary Logistic Regression

Problem Statement

The original Binary MNIST dataset is $\{(\mathbf{x}_n, \bar{y}_n)\}_{n=1}^N$ where:

$$\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ \bar{y} \in \{0, 1, \dots, 9\} \end{cases}$$

Assume we define a new problem where the labels are binarized as:

$$y_n = \begin{cases} 0 & \text{if } \bar{y}_n < 5 \\ 1 & \text{Otherwise} \end{cases}$$

Thus we have a new dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ where:

$$\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ y_n \in \{0, 1\} \end{cases}$$

Binary Logistic Regression

Model

Assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, we can use a BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b)), \quad \boldsymbol{\theta} = \{\mathbf{w}, b\}$$

Binary Logistic Regression

Model

Assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, we can use a BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b)), \quad \boldsymbol{\theta} = \{\mathbf{w}, b\}$$

where:

$\sigma(\cdot)$	Sigmoid function
\mathbf{w}	Weight vector
b	Bias value

Binary Logistic Regression

Model

Assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, we can use a BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b)), \quad \boldsymbol{\theta} = \{\mathbf{w}, b\}$$

where:

$\sigma(\cdot)$	Sigmoid function
\mathbf{w}	Weight vector
b	Bias value

Training

The model is trained easily via KL divergence. For more details please refer to [1].

Describing \mathcal{P}

Generative Modeling

Assume we just have BMINST image $\{\mathbf{x}_i\}_{i=1}^N$ without any label and we want to estimate generating distribution $p(\mathbf{x})$. Using the chain rule, we have:

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

Describing \mathcal{P}

Generative Modeling

Assume we just have BMINST image $\{\mathbf{x}_i\}_{i=1}^N$ without any label and we want to estimate generating distribution $p(\mathbf{x})$. Using the chain rule, we have:

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

Each of the factors on the right side can be seen as a BLR problem. Thus for the first term, we have:

$$p_\theta(x_1) = \text{Ber}(x_1|\sigma(b_1)), \quad b_1 \in \mathbb{R}$$

Describing \mathcal{P}

Generative Modeling

Assume we just have BMINST image $\{\mathbf{x}_i\}_{i=1}^N$ without any label and we want to estimate generating distribution $p(\mathbf{x})$. Using the chain rule, we have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_d|\mathbf{x}_{$$

Each of the factors on the right side can be seen as a BLR problem. Thus for the first term, we have:

$$p_\theta(x_1) = \text{Ber}(x_1|\sigma(b_1)), \quad b_1 \in \mathbb{R}$$

And for arbitrary d -th term with $1 < d \leq D$, we have:

$$p_\theta(x_d|\mathbf{x}_{$$

Describing \mathcal{P}

Generative Modeling

Assume we just have BMINST image $\{\mathbf{x}_i\}_{i=1}^N$ without any label and we want to estimate generating distribution $p(\mathbf{x})$. Using the chain rule, we have:

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

Each of the factors on the right side can be seen as a BLR problem. Thus for the first term, we have:

$$p_\theta(x_1) = \text{Ber}(x_1|\sigma(b_1)), b_1 \in \mathbb{R}$$

And for arbitrary d -th term with $1 < d \leq D$, we have:

$$p_\theta(x_d|\mathbf{x}_{$$

So: $\boldsymbol{\theta} = \{b_1, b_2, \mathbf{w}_2, \dots, b_D, \mathbf{w}_D\}$

Describing \mathcal{P}

Number of Parameters

For the general case of chain rule we have:

$$p(\mathbf{x}) = \overbrace{p(x_1)}^{\#1} \overbrace{p(x_2 | \mathbf{x}_{<2})}^{\#2} \dots \overbrace{p(x_D | \mathbf{x}_{$$

Describing \mathcal{P}

Number of Parameters

For the general case of chain rule we have:

$$p(\mathbf{x}) = \overbrace{p(x_1)}^{\#1} \overbrace{p(x_2 | \mathbf{x}_{<2})}^{\#2} \dots \overbrace{p(x_D | \mathbf{x}_{$$

Thus the number of parameters for \mathcal{P} is:

$$2^0 + 2^1 + 2^2 + \dots + 2^{D-1} = 2^D - 1$$

Describing \mathcal{P}_θ

Number of Parameters for \mathcal{P}_θ

For FVSBN, we have:

$$p_\theta(\mathbf{x}) = p_\theta(x_1)p_\theta(x_2|\mathbf{x}_{<2}) \dots p_\theta(x_D|\mathbf{x}_{$$

$$\Rightarrow p_\theta(\mathbf{x}) = \overbrace{\text{Ber}(x_1|\sigma(b_1))}^{\#1} \overbrace{\text{Ber}(x_2|\sigma(b_2 + w_2^T \mathbf{x}_{<2}))}^{\#2} \dots \overbrace{\text{Ber}(x_D|\sigma(b_D + \mathbf{w}_D^T \mathbf{x}_{$$

Describing \mathcal{P}_θ

Number of Parameters for \mathcal{P}_θ

For FVSBN, we have:

$$p_\theta(\mathbf{x}) = p_\theta(x_1)p_\theta(x_2|\mathbf{x}_{<2}) \dots p_\theta(x_D|\mathbf{x}_{$$

$$\Rightarrow p_\theta(\mathbf{x}) = \overbrace{\text{Ber}(x_1|\sigma(b_1))}^{\#1} \overbrace{\text{Ber}(x_2|\sigma(b_2 + w_2^T \mathbf{x}_{<2}))}^{\#2} \dots \overbrace{\text{Ber}(x_D|\sigma(b_D + w_D^T \mathbf{x}_{$$

Thus the number of parameters for \mathcal{P}_θ : $1 + 2 + 3 + \dots + D = \frac{D(D+1)}{2}$. Note that:

$$\boldsymbol{\theta} = \{b_1, b_2, \mathbf{w}_2, \dots, b_D, \mathbf{w}_D\}$$

Task: Density Estimation

\boldsymbol{x}

$$\boxed{x_1}$$

$$\boxed{x_2}$$

\vdots

$$\boxed{x_D}$$

Figure: Input D dimensional vector \boldsymbol{x}

Task: Density Estimation

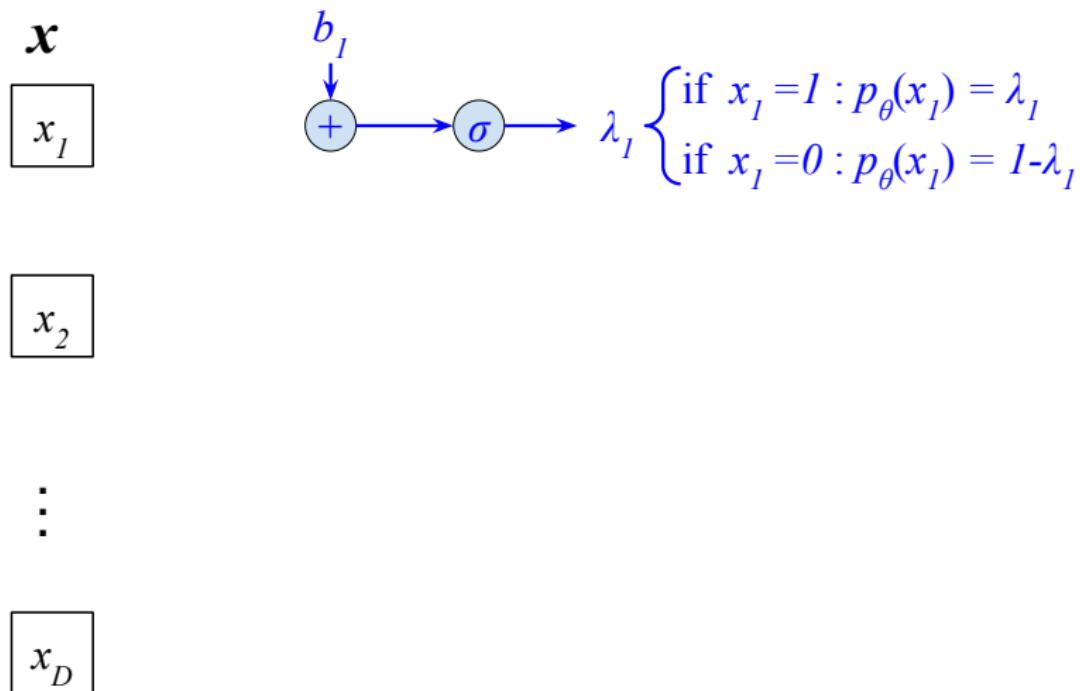


Figure: Calculations for $p_\theta(x_1)$

Task: Density Estimation

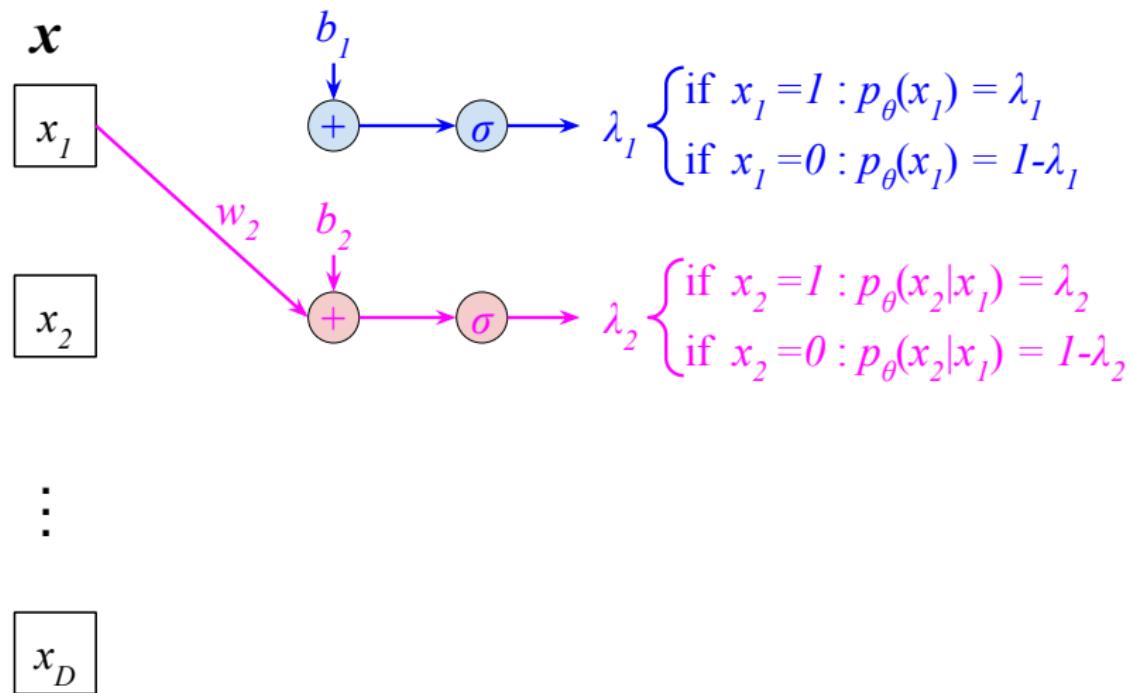


Figure: Calculations for $p_\theta(x_2|x_1)$

Task: Density Estimation

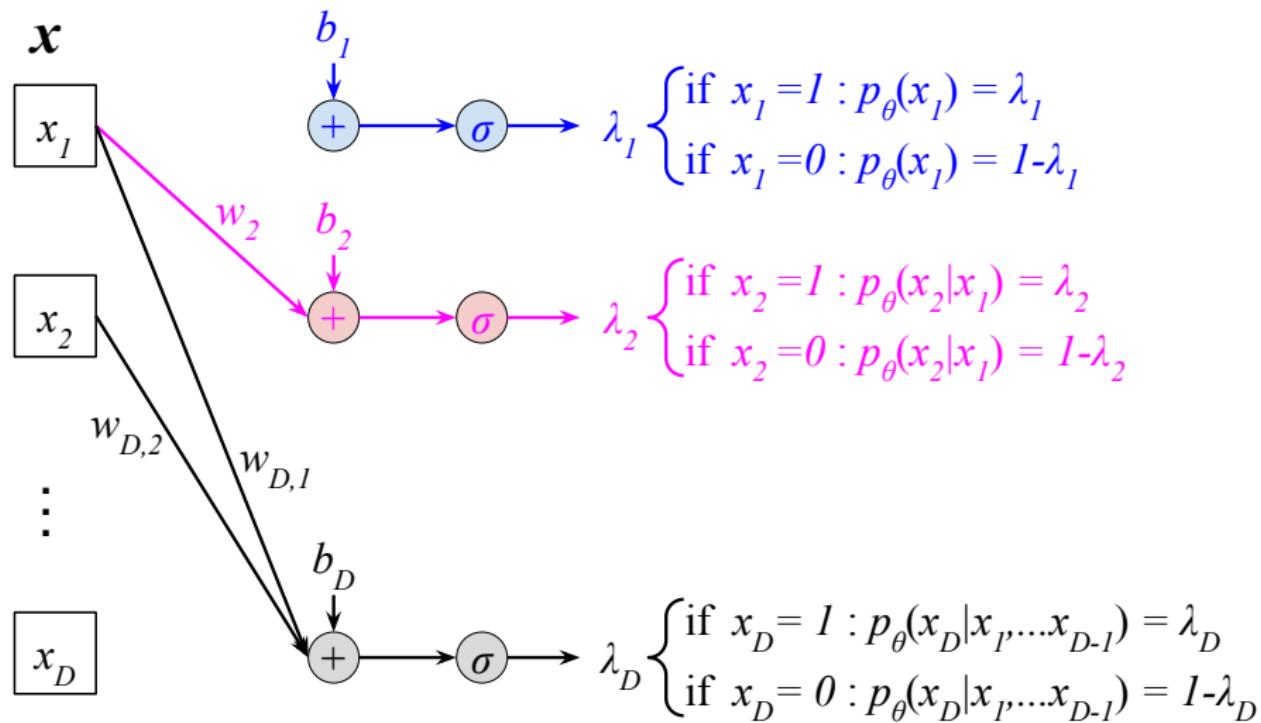


Figure: Calculations for $p_\theta(x_D | x_1 \dots x_{D-1})$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) =$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) = \overbrace{p_{\theta}(X_1 = 0)}^{\text{Ber}(X_1 | \lambda_1)}$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) = \overbrace{p_{\theta}(X_1 = 0)}^{\text{Ber}(X_1 | \lambda_1)} \overbrace{p_{\theta}(X_2 = 1 | \mathbf{x}_{<2} = [0]^T)}^{\text{Ber}(X_2 | \lambda_2(\mathbf{x}_{<2} = [0]^T))}$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) = \overbrace{p_{\theta}(X_1 = 0)}^{\text{Ber}(X_1 | \lambda_1)} \overbrace{p_{\theta}(X_2 = 1 | \mathbf{x}_{<2} = [0]^T)}^{\text{Ber}(X_2 | \lambda_2 (\mathbf{x}_{<2} = [0]^T))}$$
$$\overbrace{p_{\theta}(X_3 = 1 | \mathbf{x}_{<3} = [0, 1]^T)}^{\text{Ber}(X_3 | \lambda_3 (\mathbf{x}_{<3} = [0, 1]^T))}$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) = \overbrace{p_{\theta}(X_1 = 0)}^{\text{Ber}(X_1 | \lambda_1)} \overbrace{p_{\theta}(X_2 = 1 | \mathbf{x}_{<2} = [0]^T)}^{\text{Ber}(X_2 | \lambda_2 (\mathbf{x}_{<2} = [0]^T))}$$
$$p_{\theta}(X_3 = 1 | \mathbf{x}_{<3} = [0, 1]^T) \overbrace{p_{\theta}(X_4 = 0 | \mathbf{x}_{<4} = [0, 1, 1]^T)}^{\text{Ber}(X_4 | \lambda_4 (\mathbf{x}_{<4} = [0, 1, 1]^T))}$$

Task: Density Estimation

Density Estimation

For now, assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then for a simple case with $n = 4$ we have:

$$\begin{aligned} p_{\theta}(\mathbf{x} = [0, 1, 1, 0]^T) &= \overbrace{p_{\theta}(X_1 = 0)}^{\text{Ber}(X_1 | \lambda_1)} \overbrace{p_{\theta}(X_2 = 1 | \mathbf{x}_{<2} = [0]^T)}^{\text{Ber}(X_2 | \lambda_2 (\mathbf{x}_{<2} = [0]^T))} \\ &\quad \overbrace{p_{\theta}(X_3 = 1 | \mathbf{x}_{<3} = [0, 1]^T)}^{\text{Ber}(X_3 | \lambda_3 (\mathbf{x}_{<3} = [0, 1]^T))} \overbrace{p_{\theta}(X_4 = 0 | \mathbf{x}_{<4} = [0, 1, 1]^T)}^{\text{Ber}(X_4 | \lambda_4 (\mathbf{x}_{<4} = [0, 1, 1]^T))} \\ &= (1 - \lambda_1)(\lambda_2)(\lambda_3)(1 - \lambda_4) \end{aligned}$$

Task: Generation

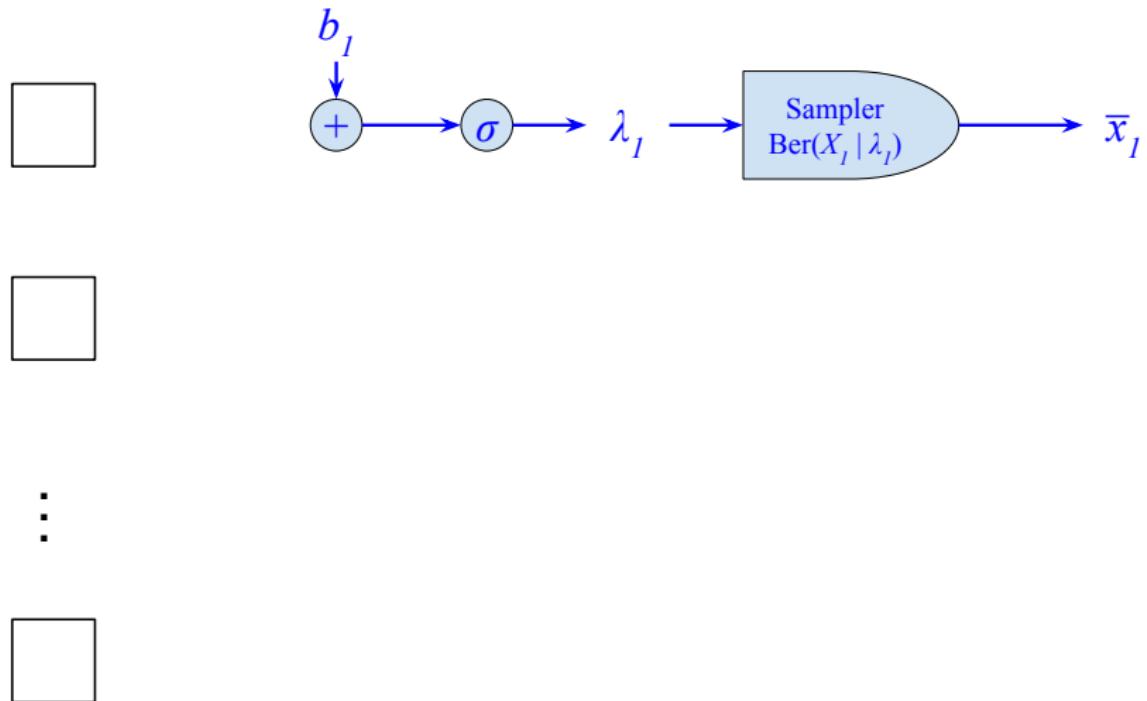


Figure: Sampling first element

Task: Generation

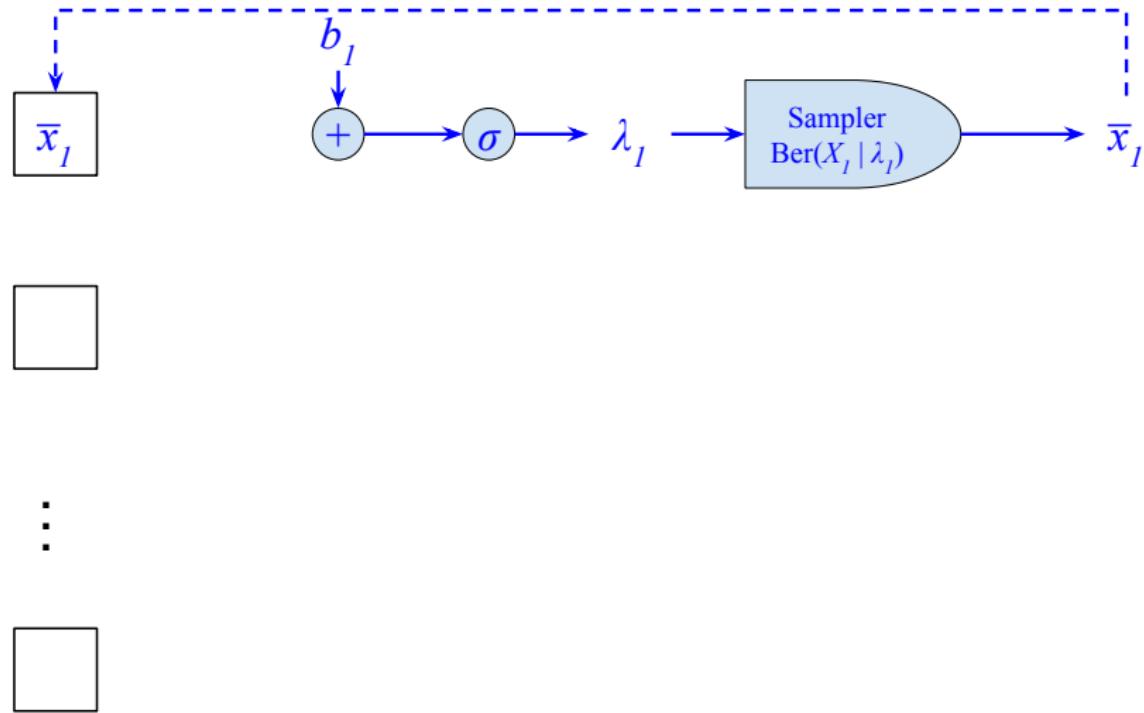


Figure: Entering the first sampled element in the vector

Task: Generation

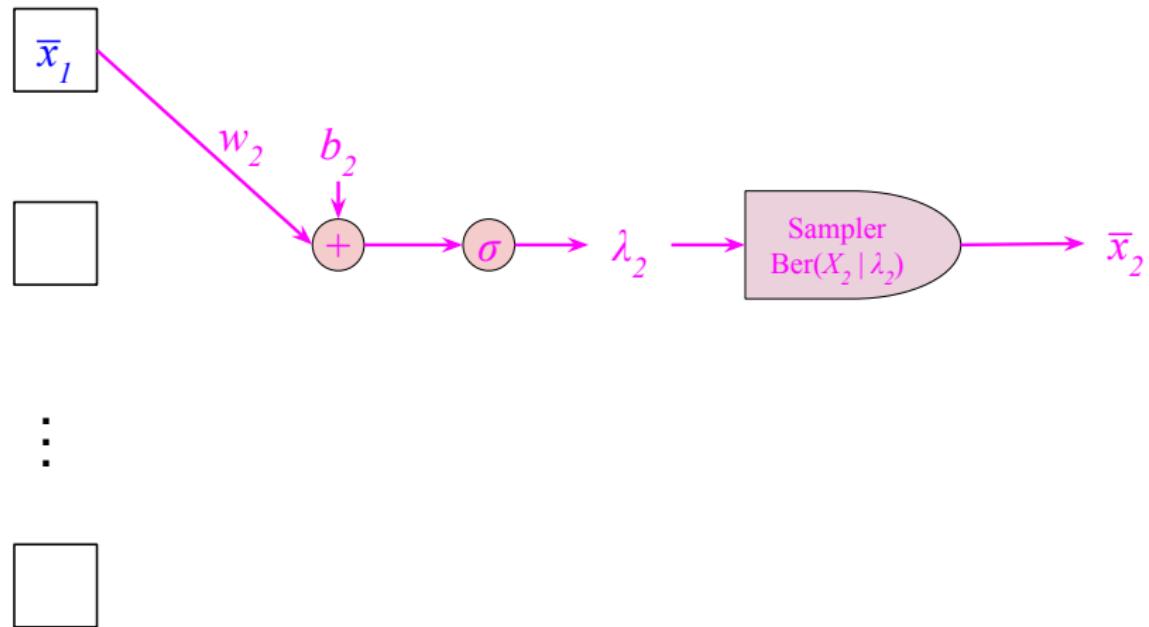


Figure: Sampling second element

Task: Generation

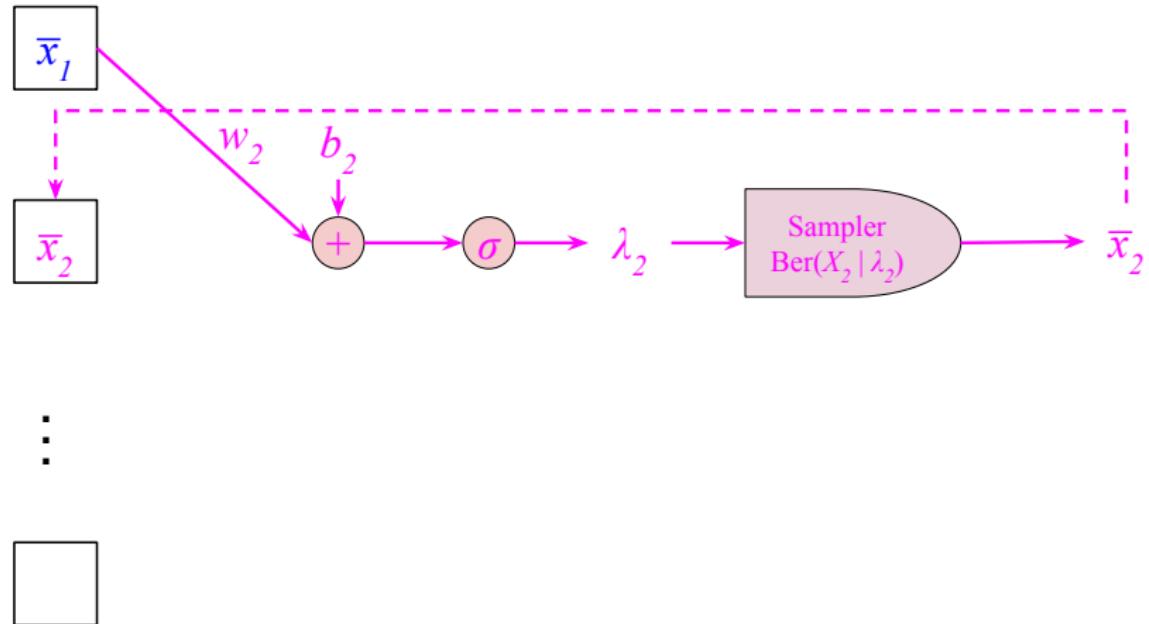


Figure: Entering the second sampled element in the vector

Task: Generation

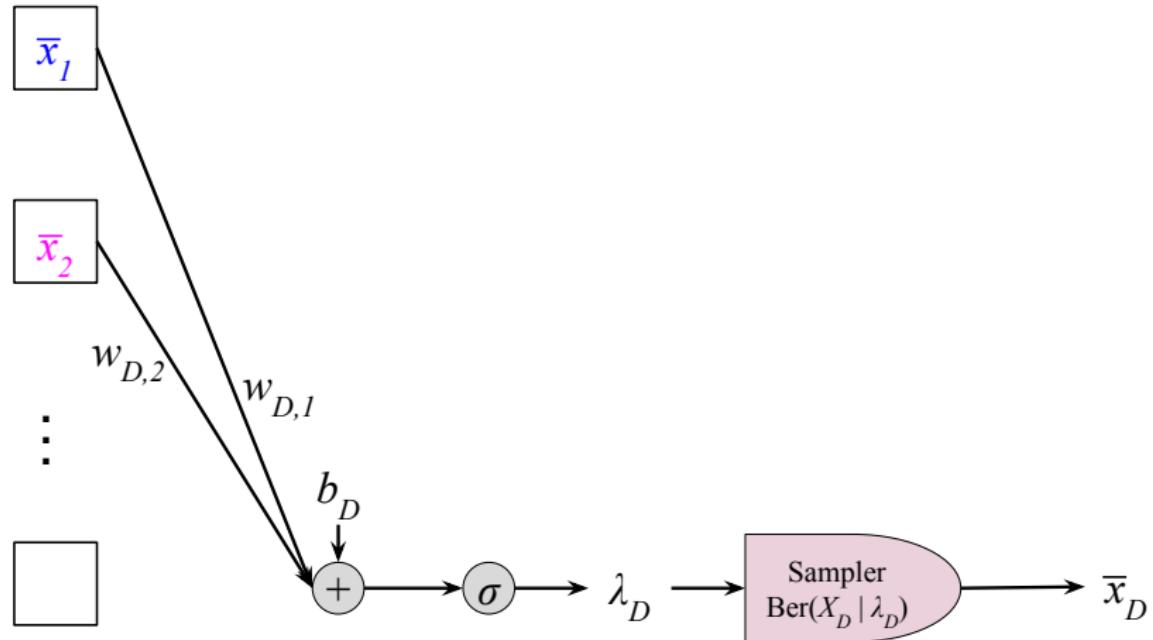


Figure: Sampling D -th element

Task: Generation

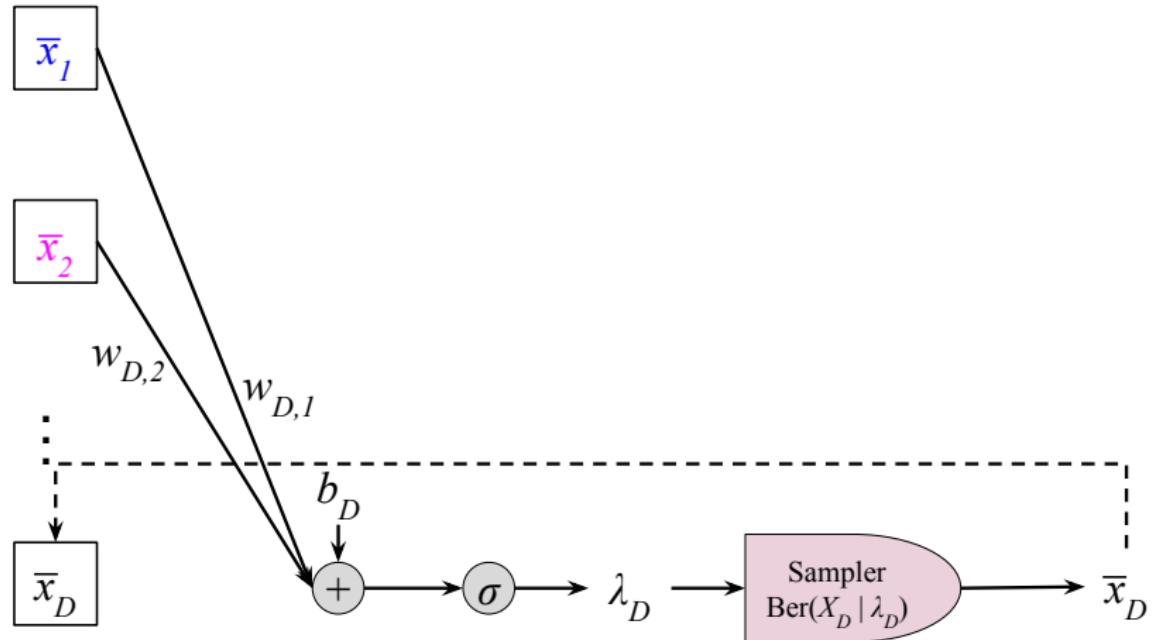


Figure: Entering the D -th sampled element in the vector

Generation

Again assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{$$

Generation

Again assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then you can sample $p(x)$ as:

$$\bar{x}_1 \sim p_{\theta}(X_1) = \text{Ber}(X_1 | \lambda_1)$$

Generation

Again assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then you can sample $p(x)$ as:

$$\bar{x}_1 \sim p_{\theta}(X_1) = \text{Ber}(X_1 | \lambda_1)$$

$$\bar{x}_2 \sim p(X_2 | \mathbf{x}_{<2} = [\bar{x}_1]^T) = \text{Ber}(X_2 | \lambda_2(\mathbf{x}_{<2} = [\bar{x}_1]^T))$$

Generation

Again assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then you can sample $p(x)$ as:

$$\bar{x}_1 \sim p_{\theta}(X_1) = \text{Ber}(X_1 | \lambda_1)$$

$$\bar{x}_2 \sim p(X_2 | \mathbf{x}_{<2} = [\bar{x}_1]^T) = \text{Ber}(X_2 | \lambda_2(\mathbf{x}_{<2} = [\bar{x}_1]^T))$$

$$\bar{x}_3 \sim p(X_3 | \mathbf{x}_{<3} = [\bar{x}_1, \bar{x}_2]^T) = \text{Ber}(X_3 | \lambda_3(\mathbf{x}_{<3} = [\bar{x}_1, \bar{x}_2]^T))$$

Generation

Again assume the parameters for FVSBN θ are given and define:

$$\lambda_d(\mathbf{x}_{)} \triangleq p_{\theta}(X_d = 1 | \mathbf{x}_{$$

Then you can sample $p(x)$ as:

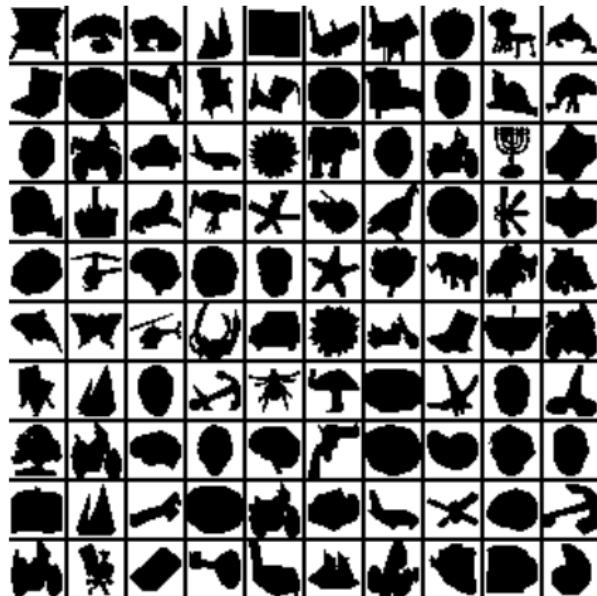
$$\bar{x}_1 \sim p_{\theta}(X_1) = \text{Ber}(X_1 | \lambda_1)$$

$$\bar{x}_2 \sim p(X_2 | \mathbf{x}_{<2} = [\bar{x}_1]^T) = \text{Ber}(X_2 | \lambda_2(\mathbf{x}_{<2} = [\bar{x}_1]^T))$$

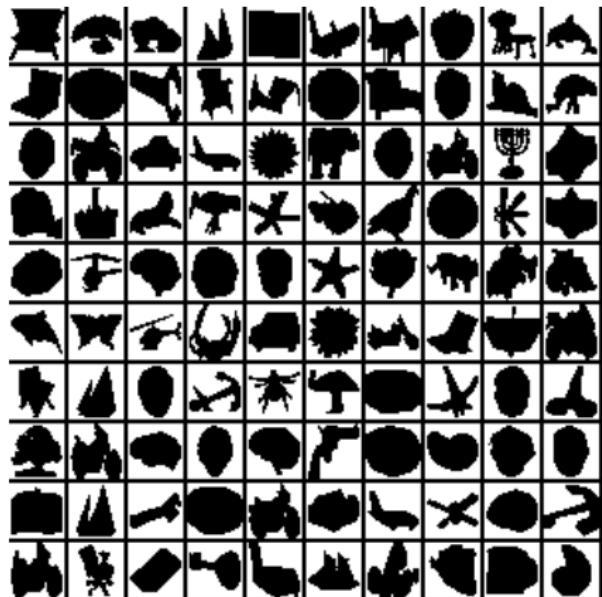
$$\bar{x}_3 \sim p(X_3 | \mathbf{x}_{<3} = [\bar{x}_1, \bar{x}_2]^T) = \text{Ber}(X_3 | \lambda_3(\mathbf{x}_{<3} = [\bar{x}_1, \bar{x}_2]^T))$$

⋮

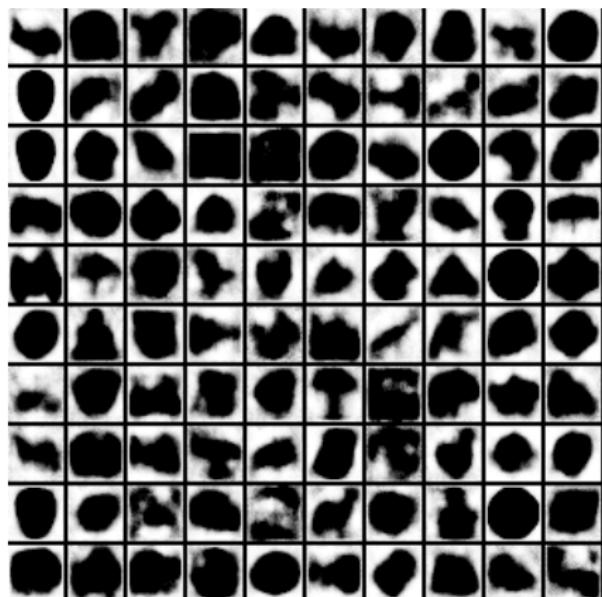
$$\bar{x}_D \sim p(X_D | \mathbf{x}_{$$



(a) Dataset samples



(a) Dataset samples



(b) Generated samples

Figure: FVSBN performance over Caltech 101 dataset (source: [2])

Section 2

Neural Autoregressive Density Estimation

BLR and Feature Transformation

BLR Limitation

Again assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ y_n \in \{0, 1\} \end{cases}$, and BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b))$$

One major limitation of the above model is

BLR and Feature Transformation

BLR Limitation

Again assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ y_n \in \{0, 1\} \end{cases}$, and BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b))$$

One major limitation of the above model is linearity with respect to input feature space, in other words:

$$p_{\theta}(y = 1|\mathbf{x}) = \frac{1}{2} \Rightarrow \mathbf{w}^T \mathbf{x} + b = 0$$

To improve expressivity, we can use

BLR and Feature Transformation

BLR Limitation

Again assume dataset $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$, $\begin{cases} \mathbf{x}_n \in \{0, 1\}^{28 \times 28} \\ y_n \in \{0, 1\} \end{cases}$, and BLR model as:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{x} + b))$$

One major limitation of the above model is linearity with respect to input feature space, in other words:

$$p_{\theta}(y = 1|\mathbf{x}) = \frac{1}{2} \Rightarrow \mathbf{w}^T \mathbf{x} + b = 0$$

To improve expressivity, we can use feature transformation:

$$p_{\theta}(y|\mathbf{x}) = \text{Ber}(y|\sigma(\mathbf{w}^T \mathbf{h}(\mathbf{x}) + b))$$

where $\mathbf{h}(\cdot)$ is the feature transformation function.

BLR and Feature Transformation

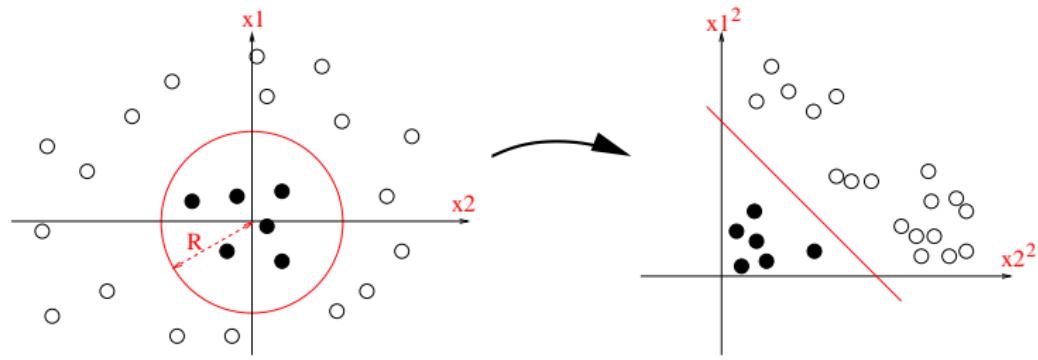


Figure: Original feature space $[x_1, x_2]^T$ (left) and transformed feature space $[x_1^2, x_2^2]^T$ (right) (source: [3]). Note that the decision boundary (red line) in the transformed space is the nonlinear decision boundary (red line) in the original space.

Generative Modeling

Based on the chain rule, we have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_d|\mathbf{x}_{)} \dots p(x_D|\mathbf{x}_{)$$

Generative Modeling

Based on the chain rule, we have:

$$p(\mathbf{x}) = p(x_1)p(x_2|\mathbf{x}_{<2}) \dots p(x_d|\mathbf{x}_{)} \dots p(x_D|\mathbf{x}_{$$

Now assume transformed feature \mathbf{h}_d defined as:

$$\mathbf{h}_d = \sigma(\mathbf{W}_d \mathbf{x}_{} + \mathbf{c}_d), \begin{cases} \mathbf{W}_d \in \mathbb{R}^{D_h \times (d-1)} \\ \mathbf{c}_d, \mathbf{h}_d \in \mathbb{R}^{D_h} \\ \sigma(\cdot) \text{ is arbitrary element-wise nonlinearity} \end{cases}$$

Generative Modeling

Thus for the first term, we have:

$$p_{\theta}(x_1) = \text{Ber}(x_1 | \sigma(b_1)), \quad b_1 \in \mathbb{R}$$

Generative Modeling

Thus for the first term, we have:

$$p_{\theta}(x_1) = \text{Ber}(x_1 | \sigma(b_1)), \quad b_1 \in \mathbb{R}$$

And for arbitrary d -th term with $1 < d \leq D$, we have:

$$p_{\theta}(x_d | \mathbf{x}_{)} = \text{Ber}(x_i | \sigma(\mathbf{v}_d^T \mathbf{h}_d + b_d))$$

Task: Density Estimation

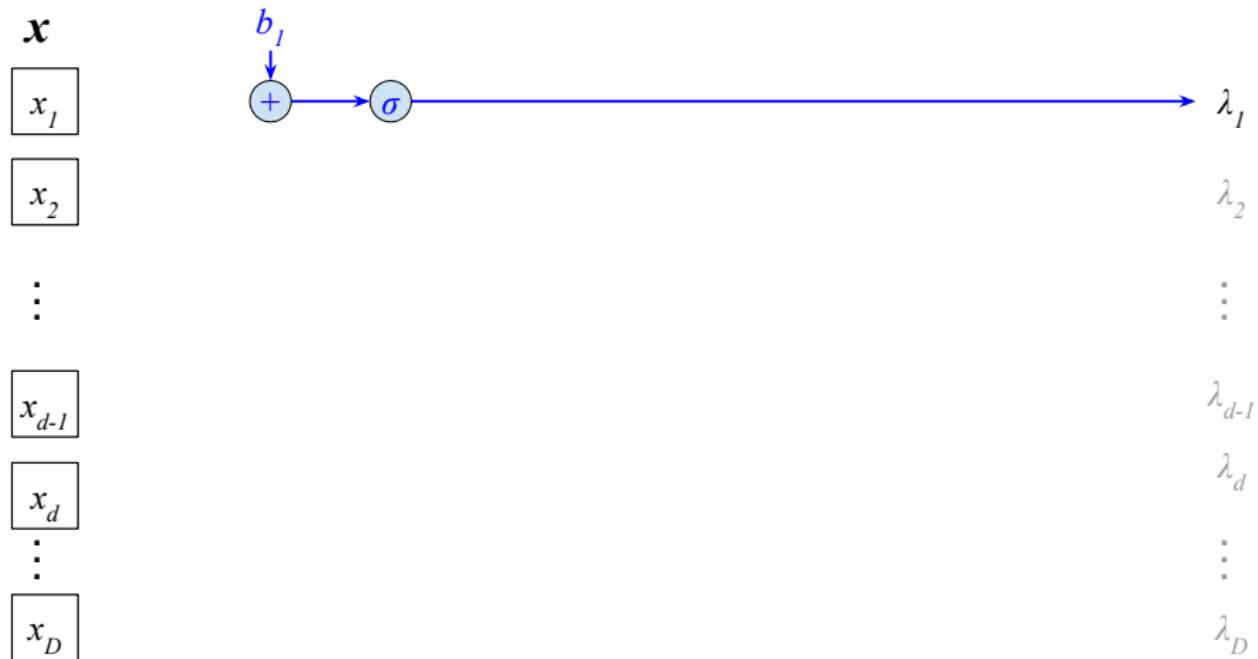


Figure: Calculations for $p_\theta(x_1)$

Task: Density Estimation

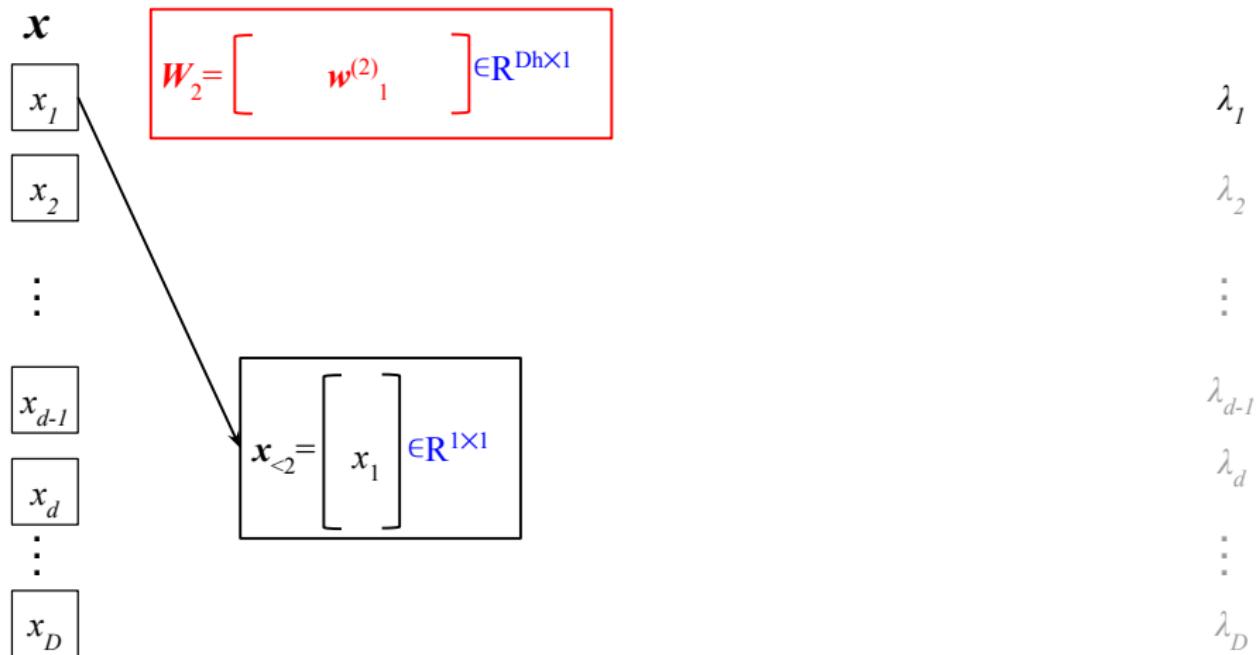


Figure: Calculations for $p_\theta(x_2|x_1)$

Task: Density Estimation

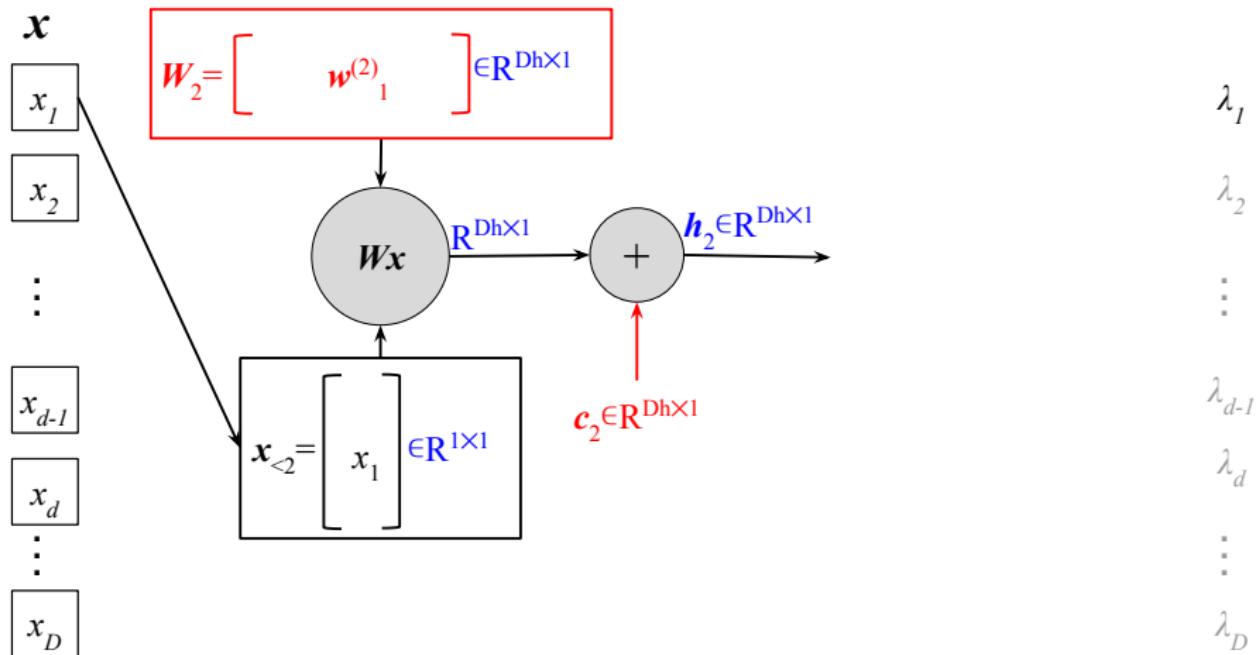


Figure: Calculations for $p_\theta(x_2|x_1)$

Task: Density Estimation

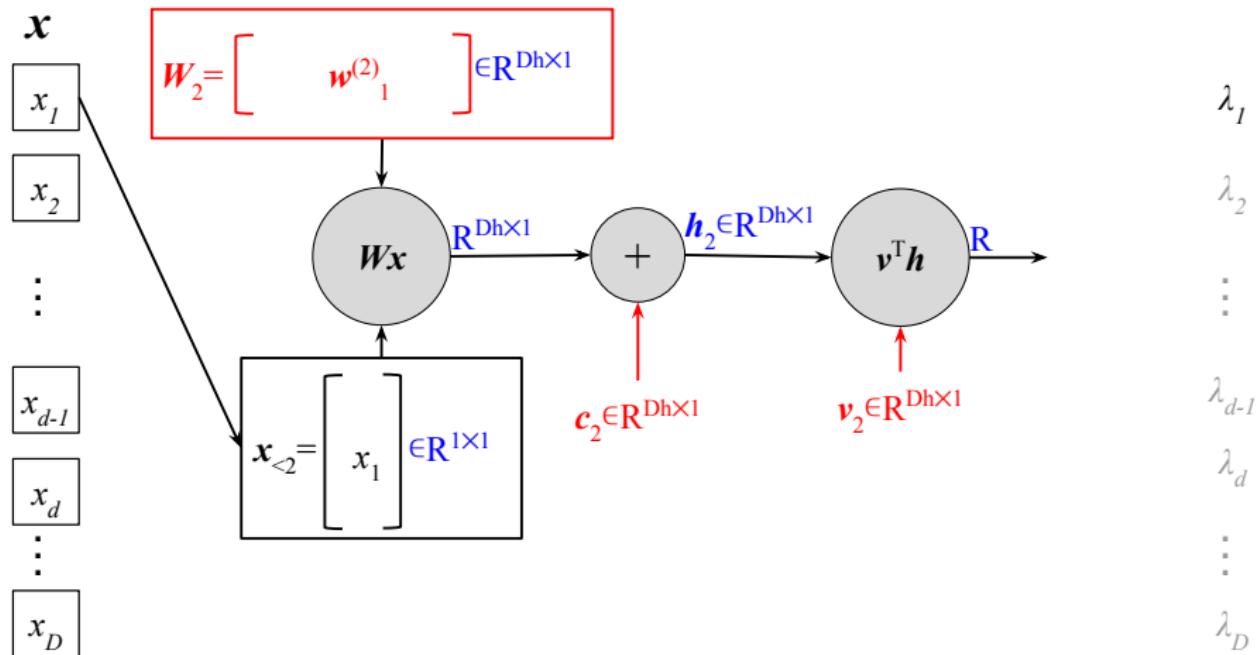
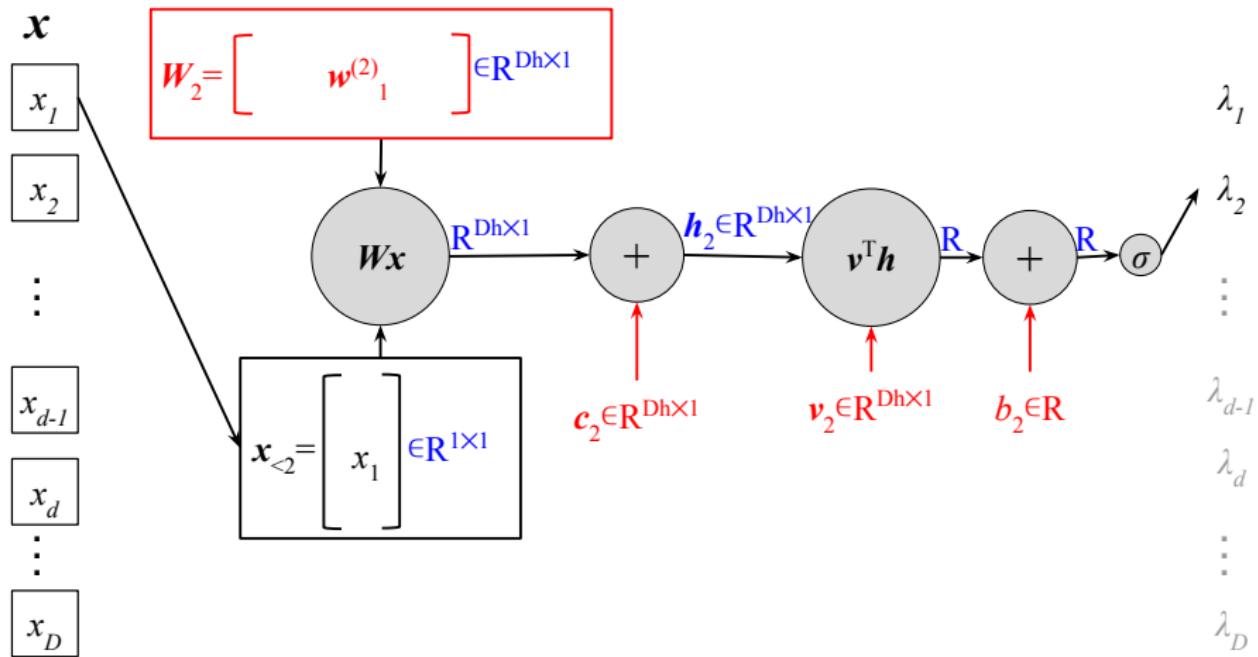


Figure: Calculations for $p_\theta(x_2|x_1)$

Task: Density Estimation



Task: Density Estimation

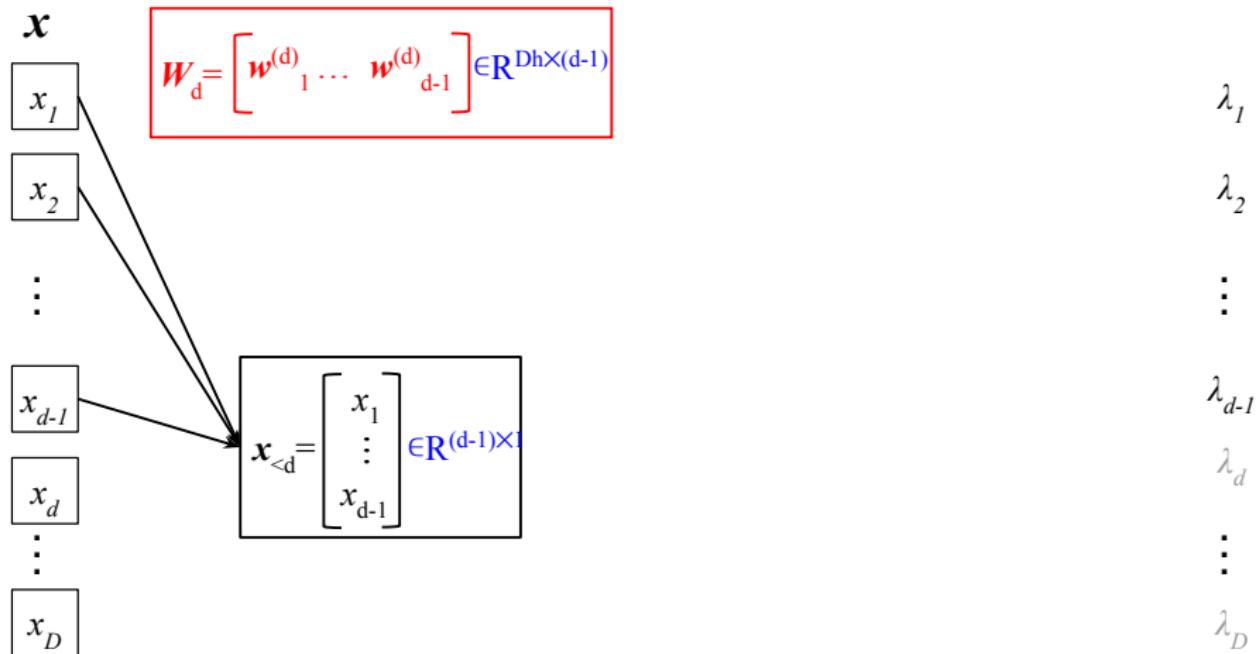


Figure: Calculations for $p_\theta(x_d | x_1, \dots, x_{d-1})$

Task: Density Estimation

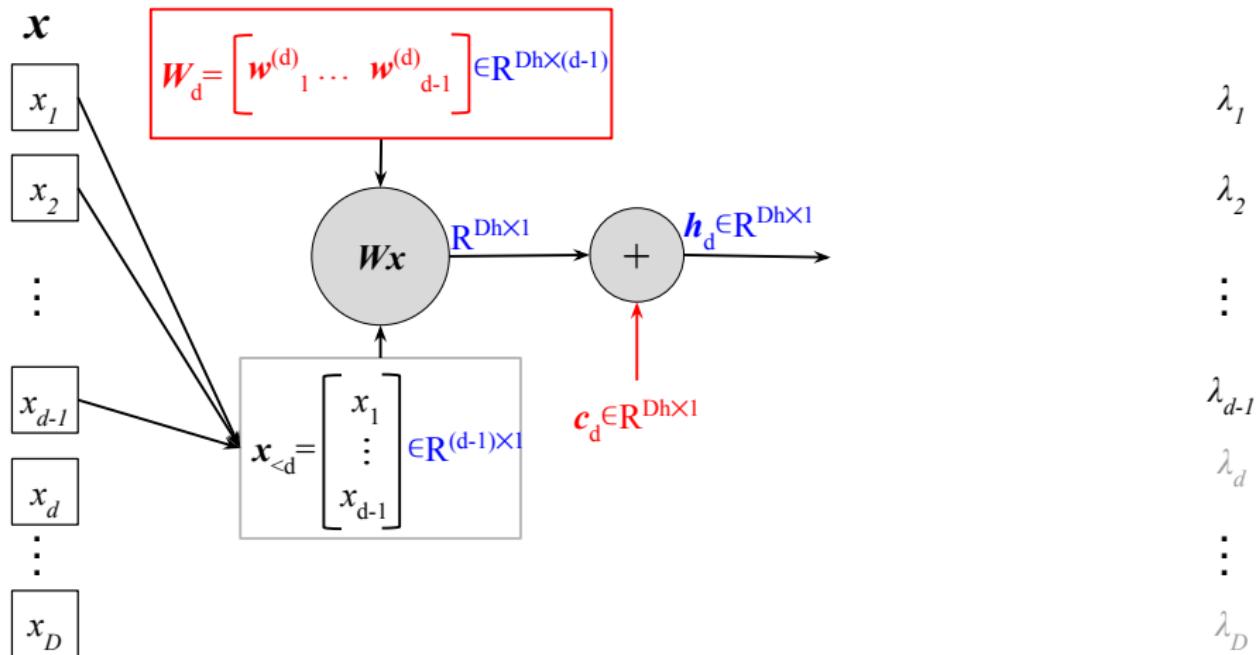


Figure: Calculations for $p_\theta(x_d|x_1, \dots, x_{d-1})$

Task: Density Estimation

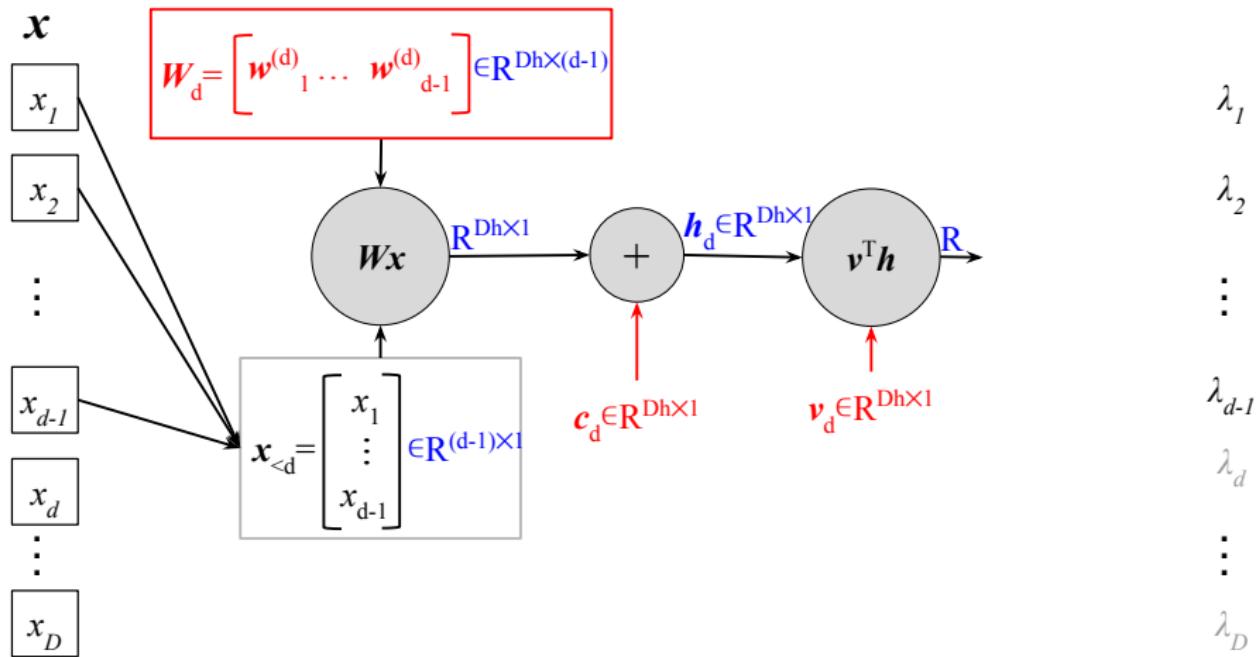


Figure: Calculations for $p_\theta(x_d | x_1, \dots, x_{d-1})$

Task: Density Estimation

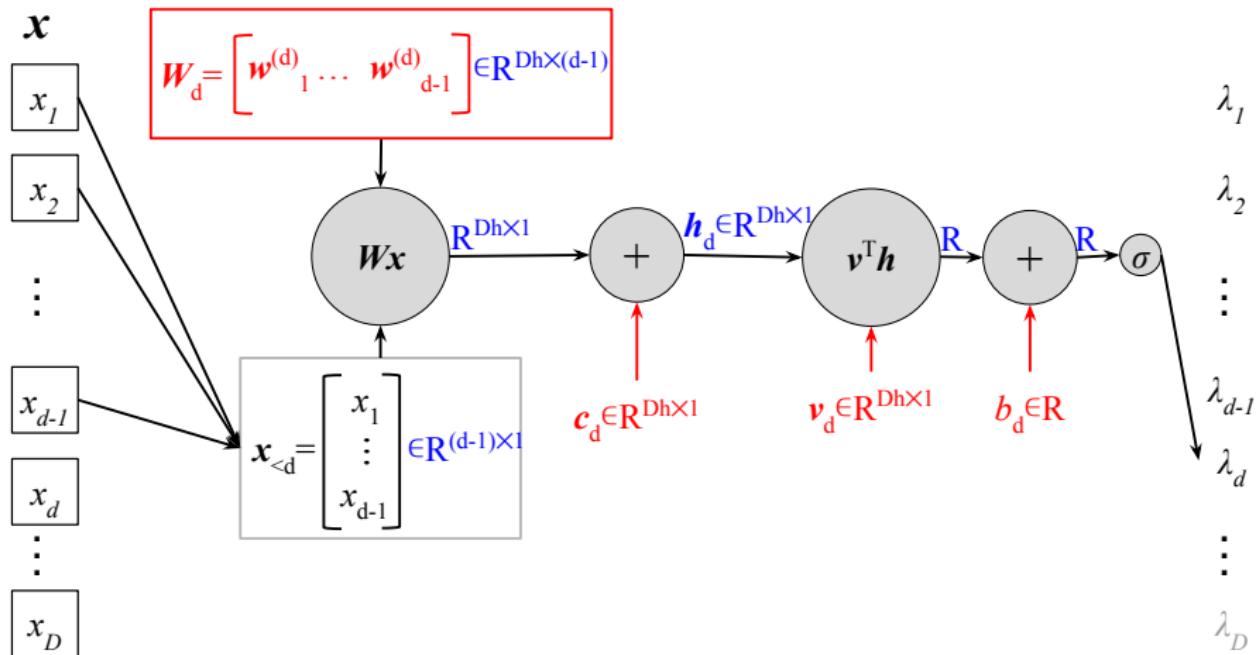


Figure: Calculations for $p_\theta(x_d|x_1, \dots, x_{d-1})$

Task: Density Estimation

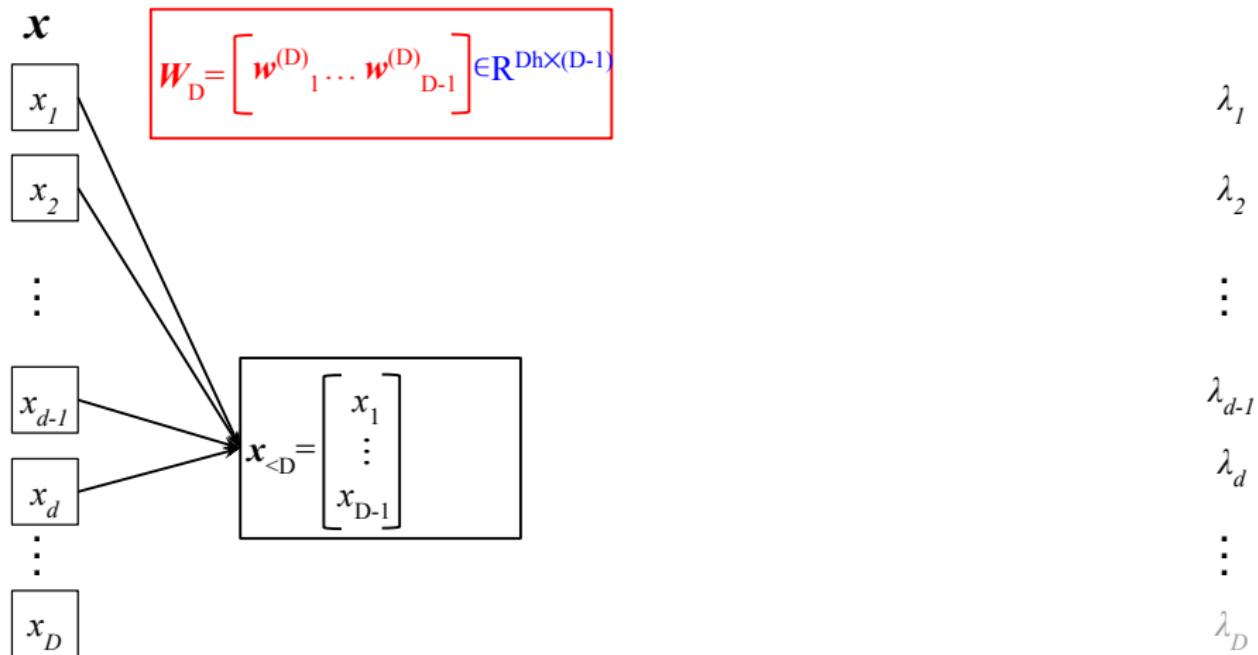


Figure: Calculations for $p_\theta(x_D | x_1, \dots, x_{D-1})$

Task: Density Estimation

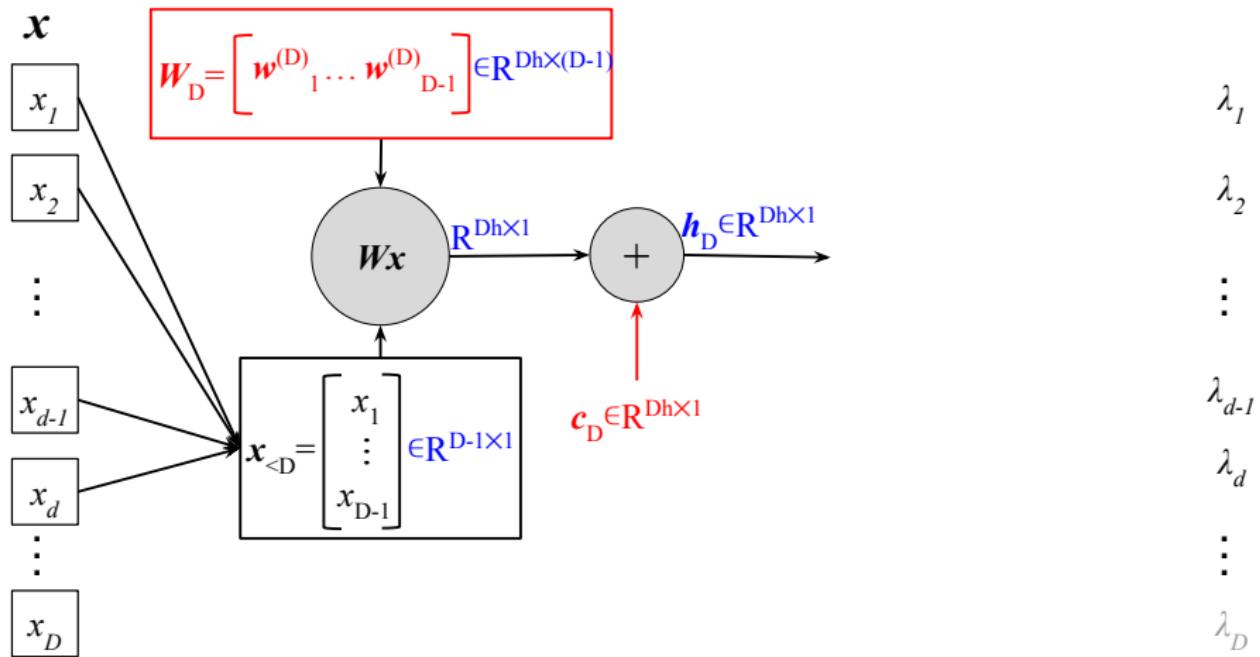


Figure: Calculations for $p_\theta(x_D | x_1, \dots, x_{D-1})$

Task: Density Estimation

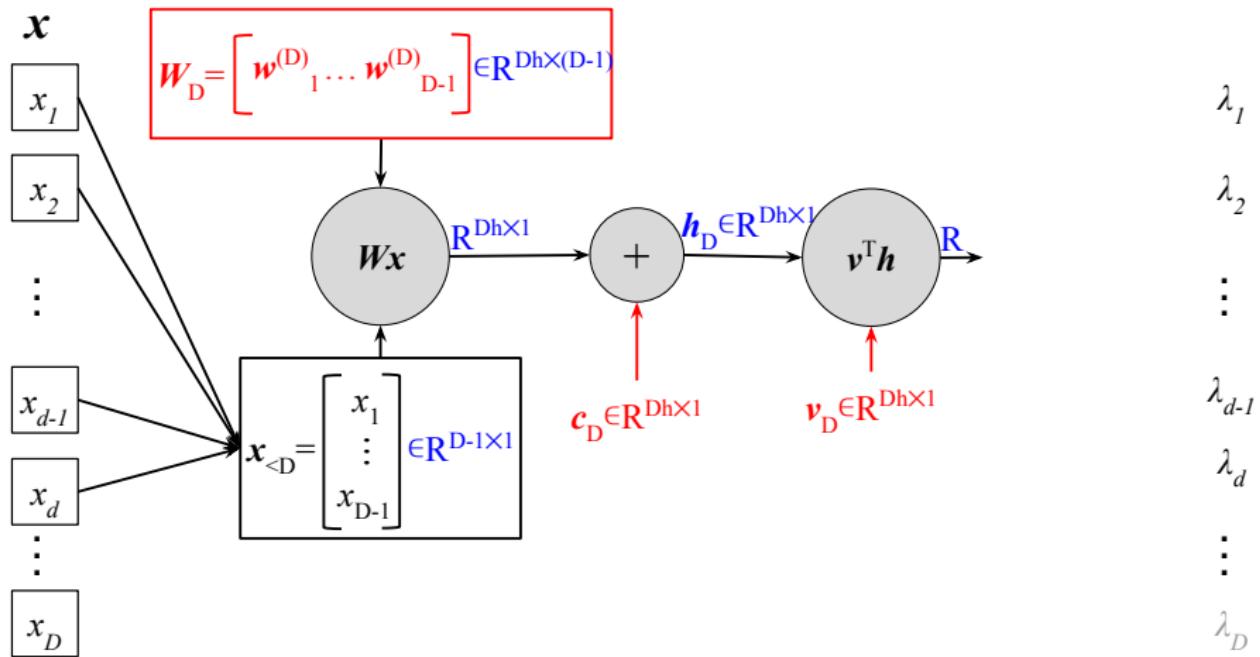


Figure: Calculations for $p_\theta(x_D | x_1, \dots, x_{D-1})$

Task: Density Estimation

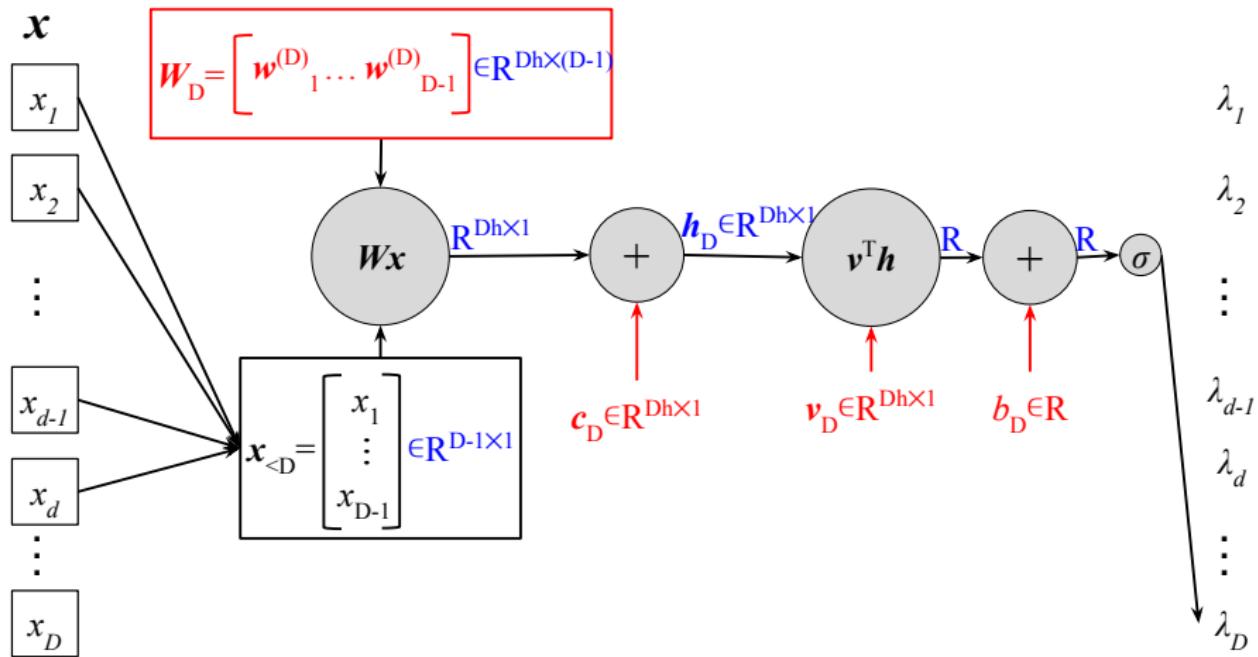


Figure: Calculations for $p_\theta(x_D | x_1, \dots, x_{D-1})$

Weight Sharing

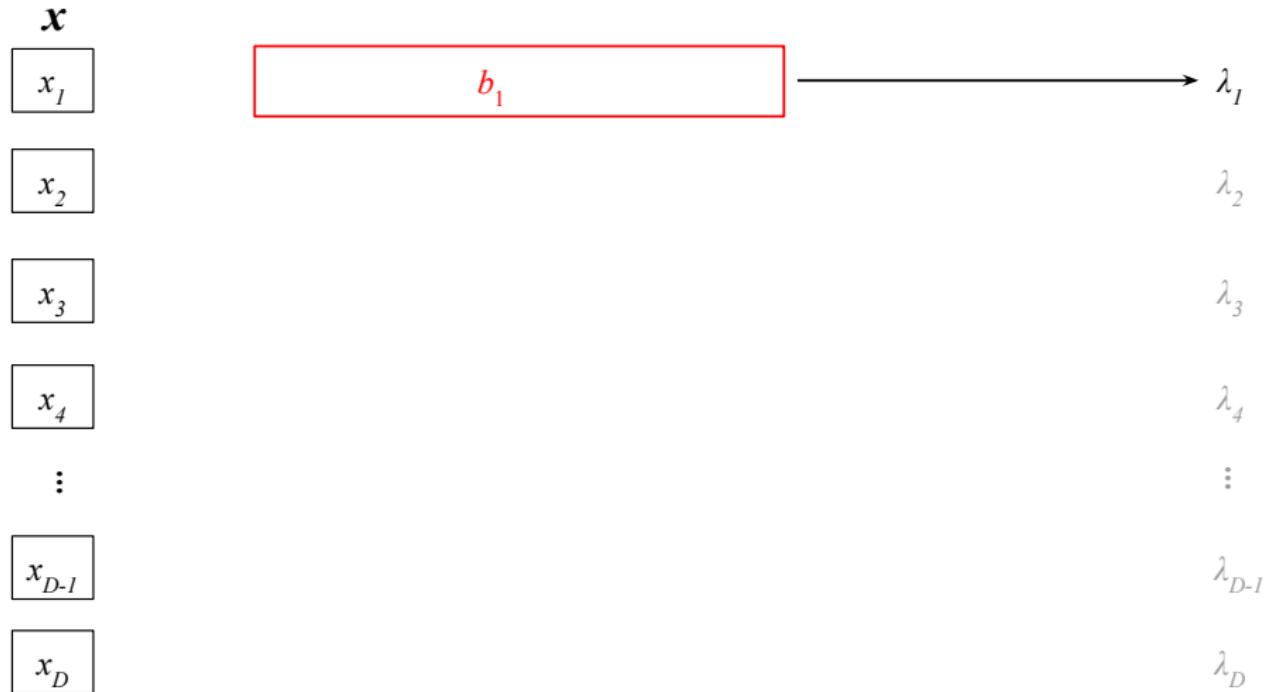


Figure: Calculations for $p_\theta(x_1)$

Weight Sharing

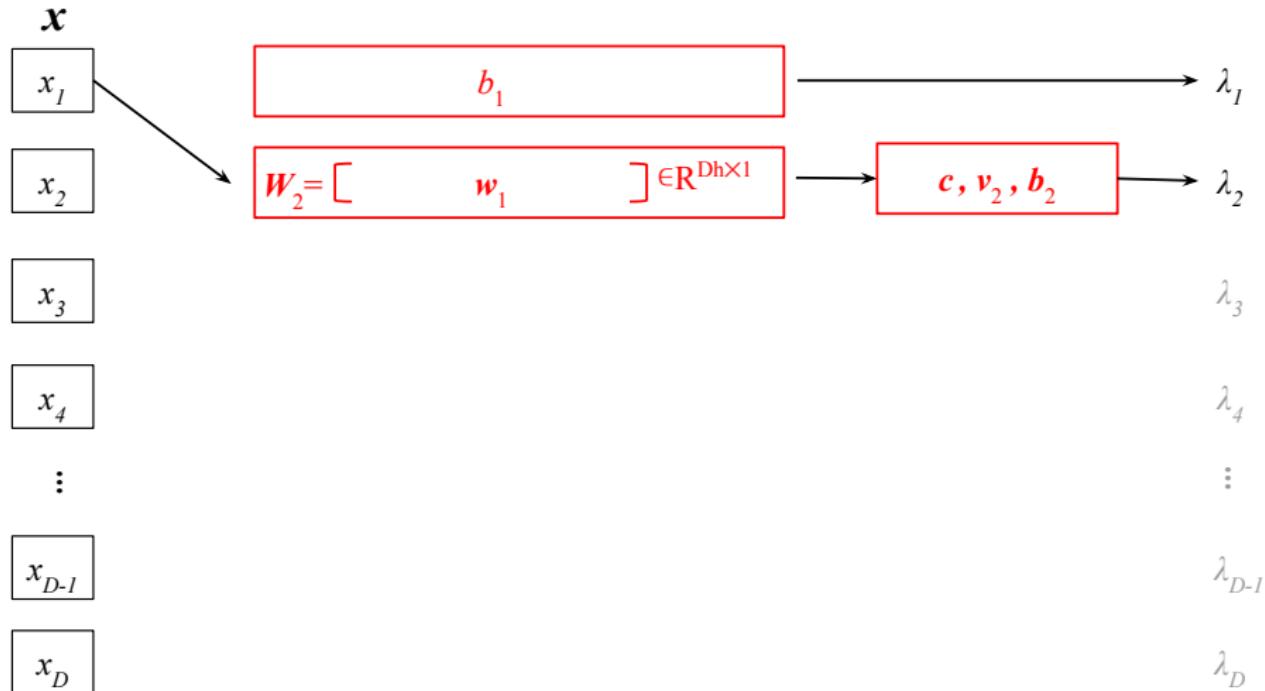


Figure: Calculations for $p_\theta(x_2|x_1)$

Weight Sharing

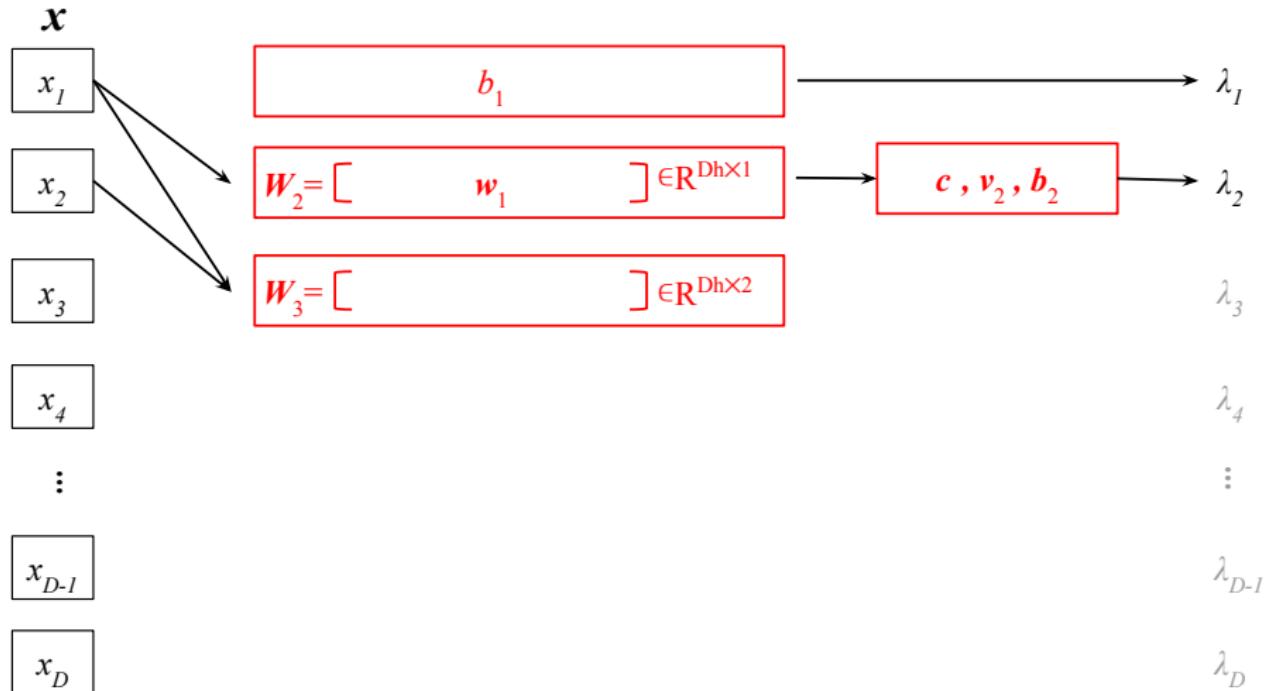


Figure: Calculations for $p_\theta(x_3|x_1, x_2)$

Weight Sharing

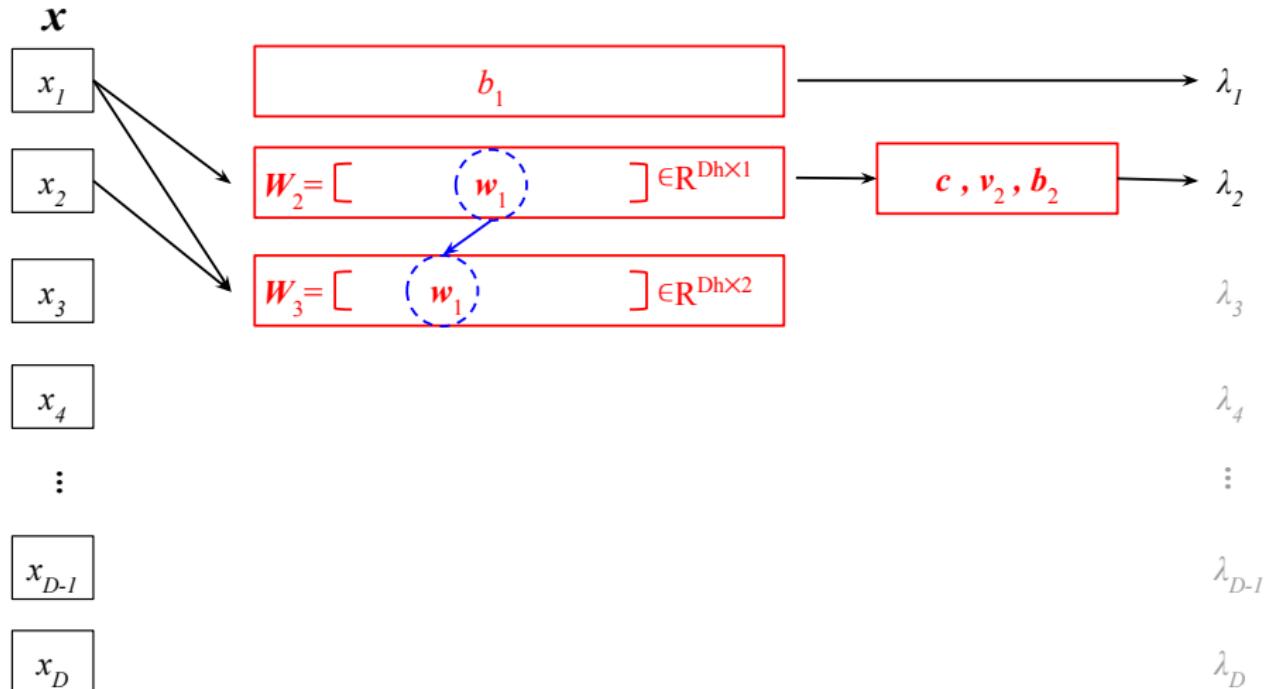


Figure: Calculations for $p_\theta(x_3|x_1, x_2)$

Weight Sharing

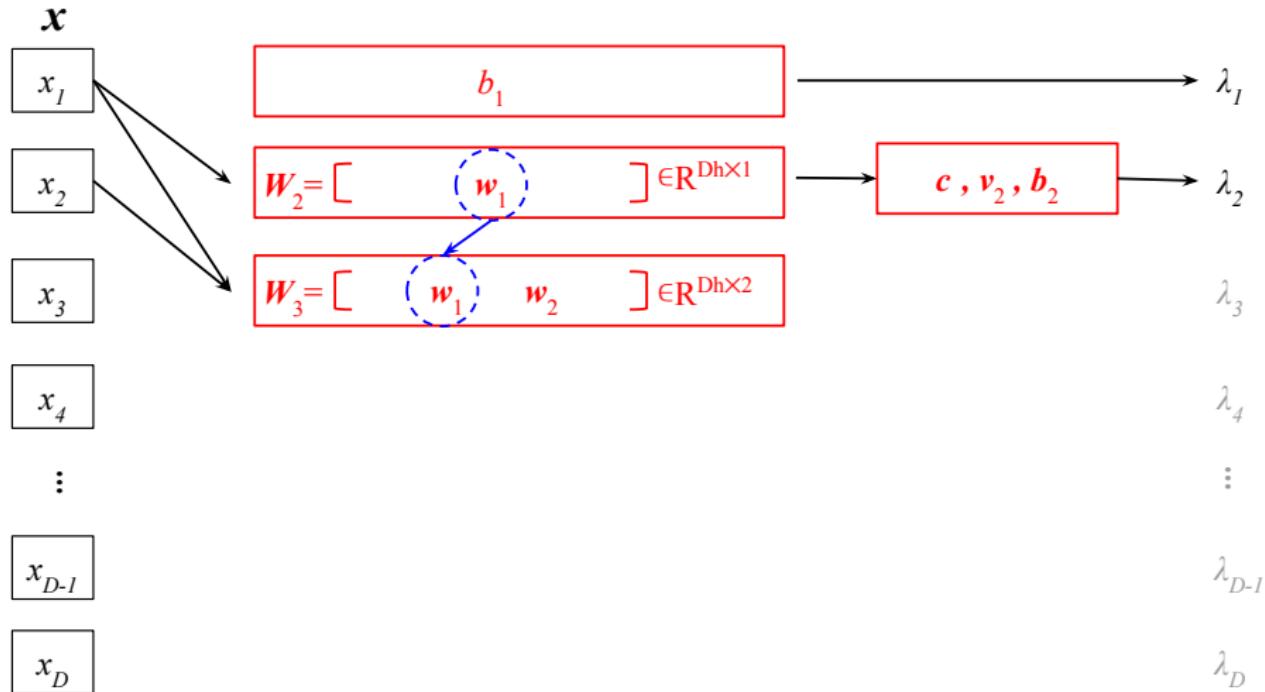


Figure: Calculations for $p_\theta(x_3|x_1, x_2)$

Weight Sharing

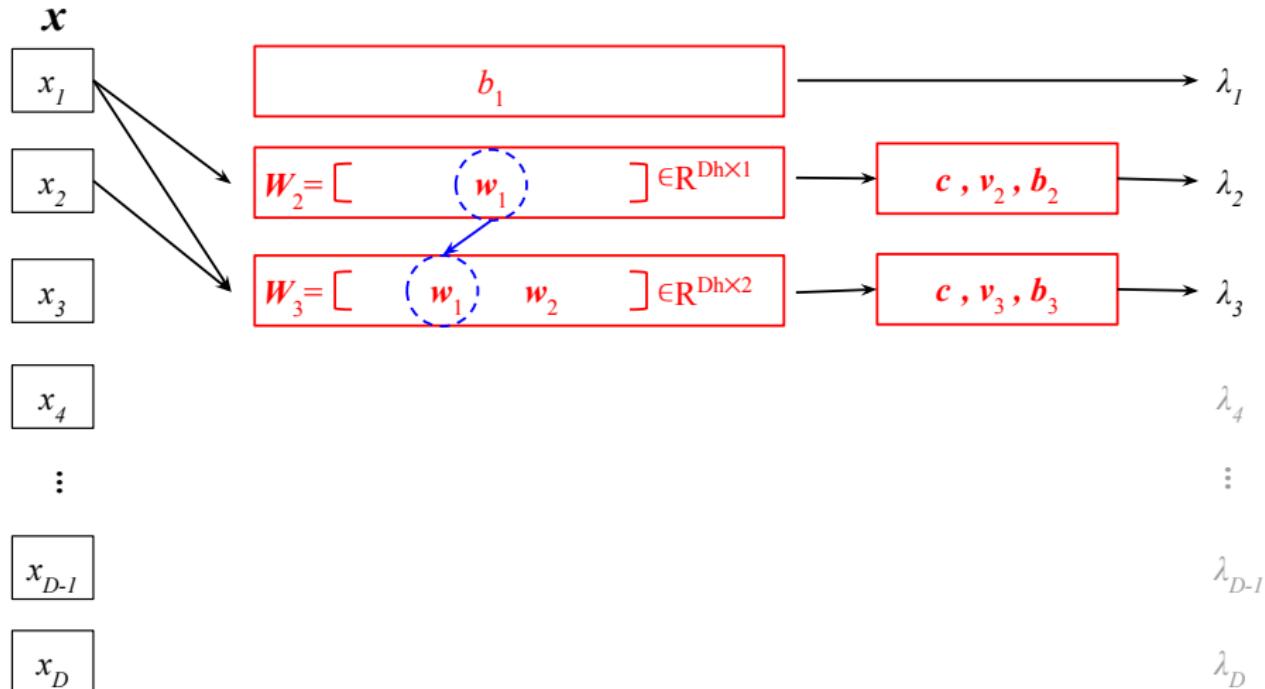


Figure: Calculations for $p_\theta(x_3|x_1, x_2)$

Weight Sharing

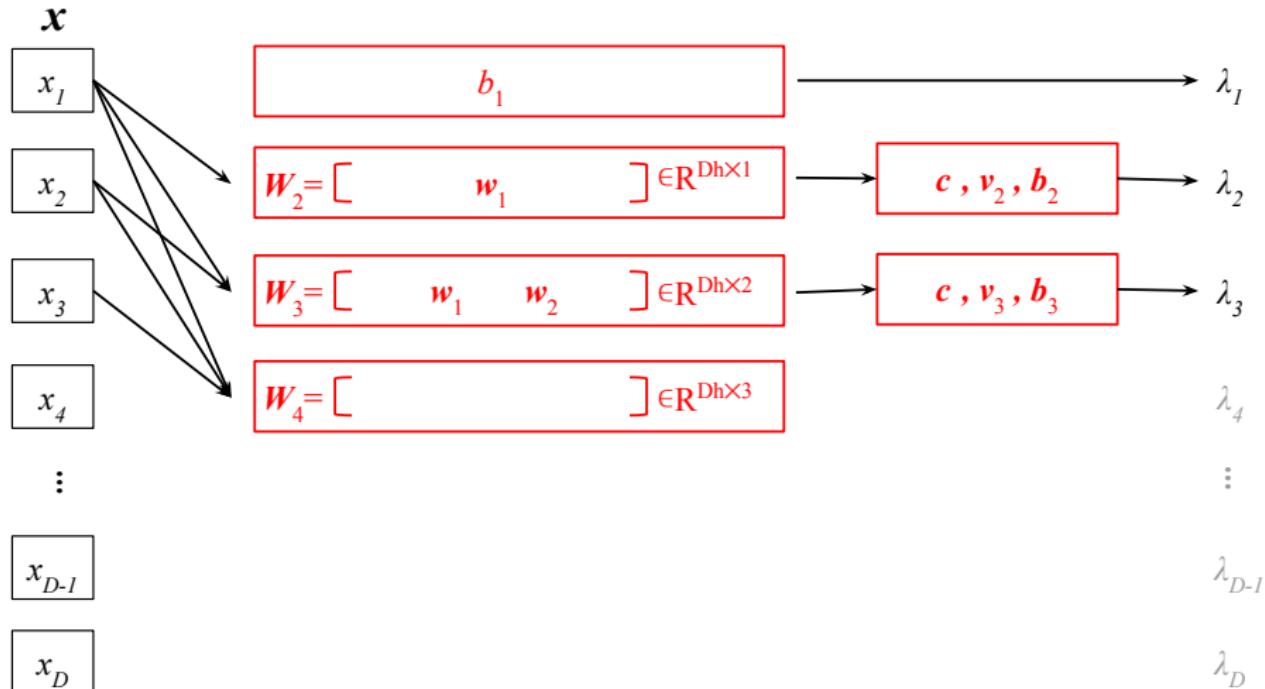


Figure: Calculations for $p_\theta(x_4|x_1, x_2, x_3)$

Weight Sharing

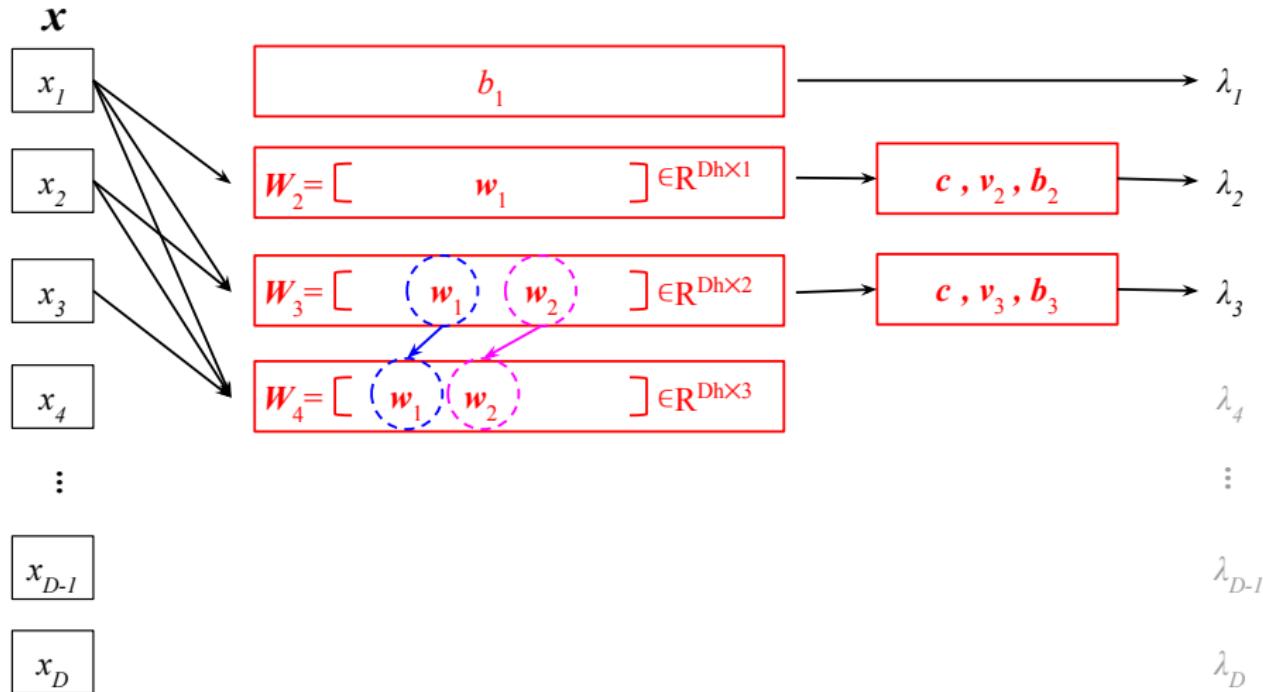


Figure: Calculations for $p_\theta(x_4|x_1, x_2, x_3)$

Weight Sharing

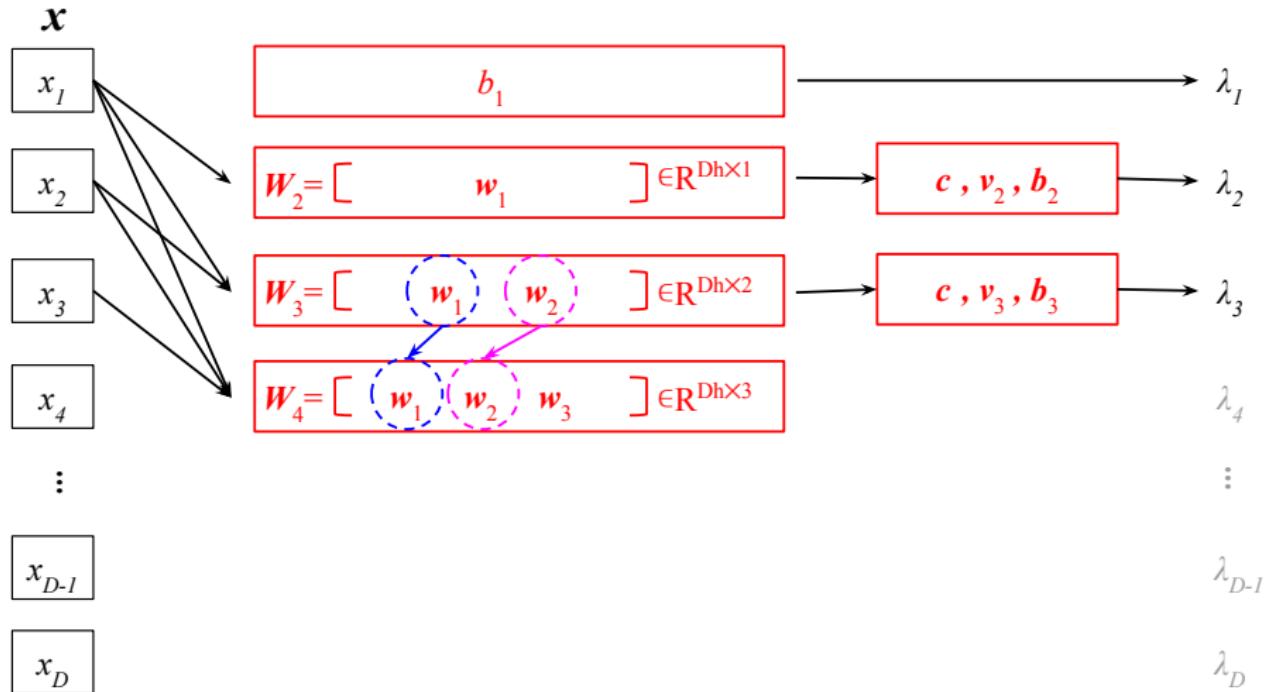


Figure: Calculations for $p_\theta(x_4|x_1, x_2, x_3)$

Weight Sharing

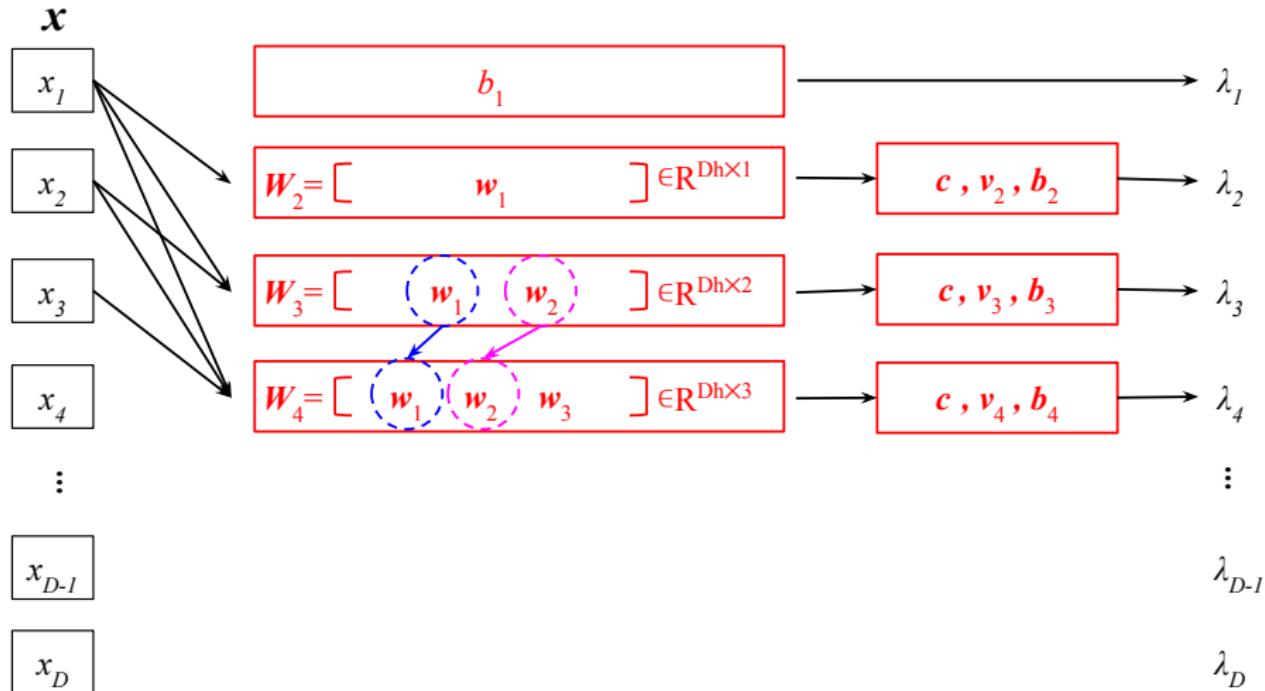


Figure: Calculations for $p_\theta(x_4|x_1, x_2, x_3)$

Weight Sharing

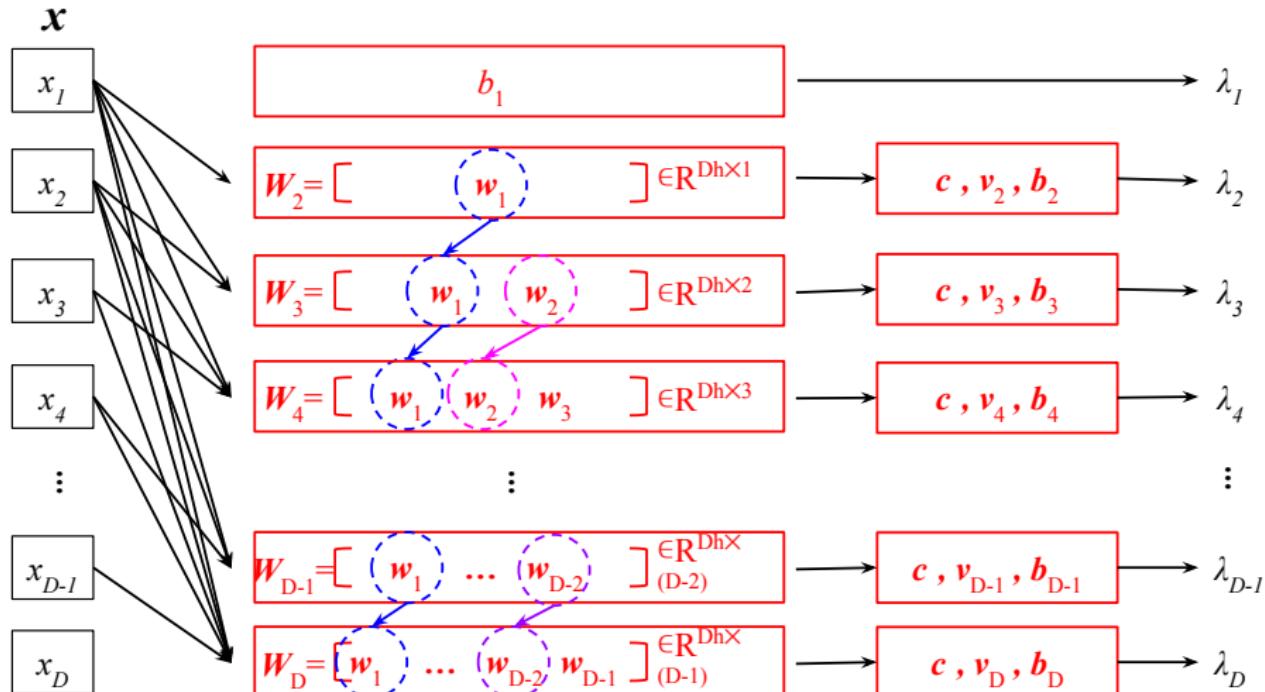


Figure: Calculations for $p_\theta(\mathbf{x})$

Weight Sharing

To reduce the number of parameters and speed up computations, the weights for generating hidden representation are shared as follows:

Weight Sharing

To reduce the number of parameters and speed up computations, the weights for generating hidden representation are shared as follows:

$$\mathbf{h}_2 = \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_2 \\ \mathbf{w}_1 \end{bmatrix}}_{\text{Weights}} \underbrace{\begin{bmatrix} \mathbf{x}_{<2} \\ \mathbf{x}_1 \end{bmatrix}}_{\text{Inputs}} + \mathbf{c} \right)$$

Weight Sharing

To reduce the number of parameters and speed up computations, the weights for generating hidden representation are shared as follows:

$$\mathbf{h}_2 = \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}}_{\text{shared weights}} \underbrace{\begin{bmatrix} x_1 \\ x_{<2} \end{bmatrix}}_{\text{input}} + \mathbf{c} \right), \mathbf{h}_3 = \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 \\ \mathbf{w}_3 \end{bmatrix}}_{\text{shared weights}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_{<3} \end{bmatrix}}_{\text{input}} + \mathbf{c} \right)$$

Weight Sharing

To reduce the number of parameters and speed up computations, the weights for generating hidden representation are shared as follows:

$$\begin{aligned}\mathbf{h}_2 &= \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}}_{\mathbf{W}_2} \underbrace{\begin{bmatrix} x_1 \\ x_{<2} \end{bmatrix}}_{\mathbf{x}_{<2}} + \mathbf{c} \right), \quad \mathbf{h}_3 = \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 \\ \mathbf{w}_3 & \mathbf{w}_4 \end{bmatrix}}_{\mathbf{W}_3} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_{<3} \end{bmatrix}}_{\mathbf{x}_{<3}} + \mathbf{c} \right) \\ \mathbf{h}_D &= \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_{D-1} \end{bmatrix}}_{\mathbf{W}_D} \underbrace{\begin{bmatrix} x_1 \\ \dots \\ x_{D-1} \end{bmatrix}}_{\mathbf{x}_{<D}} + \mathbf{c} \right)\end{aligned}$$

Weight Sharing

To reduce the number of parameters and speed up computations, the weights for generating hidden representation are shared as follows:

$$\begin{aligned} \mathbf{h}_2 &= \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}}_{\mathbf{W}_2} \underbrace{\begin{bmatrix} x_1 \\ x_{<2} \end{bmatrix}}_{\mathbf{x}_{<2}} + \mathbf{c} \right), \quad \mathbf{h}_3 = \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 & \mathbf{w}_2 \\ \mathbf{w}_3 \end{bmatrix}}_{\mathbf{W}_3} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_{<3} \end{bmatrix}}_{\mathbf{x}_{<3}} + \mathbf{c} \right) \\ \mathbf{h}_D &= \sigma \left(\underbrace{\begin{bmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_{D-1} \end{bmatrix}}_{\mathbf{W}_D} \underbrace{\begin{bmatrix} x_1 \\ \vdots \\ x_{D-1} \end{bmatrix}}_{\mathbf{x}_{<D}} + \mathbf{c} \right) \end{aligned}$$

Thus all the weights needed to calculate $\mathbf{h}_2, \dots, \mathbf{h}_D$ are available in \mathbf{W}_D which we represent by \mathbf{W} and also $\mathbf{c}_1 = \mathbf{c}_2 = \dots = \mathbf{c}_D = \mathbf{c}$. So we can write:

$$\mathbf{h}_d = \sigma(\mathbf{W}_{:<d} \mathbf{x}_{<d} + \mathbf{c})$$

Number of Parameters

The total number of parameters in NADE is :

Number of Parameters

The total number of parameters in NADE is :

$$\left. \begin{array}{l} \mathbf{W} \in \mathbb{R}^{D_h \times (D-1)} \\ \mathbf{c} \in \mathbb{R}^{D_h} \end{array} \right\} \Rightarrow D \times D_h$$
$$\left. \begin{array}{l} \{\mathbf{v}_d \in \mathbb{R}^{D_h}\}_{d=2}^D \\ \{b_d \in \mathbb{R}\}_{d=1}^D \end{array} \right\} \Rightarrow (D-1)D_h + D$$

Number of Parameters

The total number of parameters in NADE is :

$$\left. \begin{array}{l} \mathbf{W} \in \mathbb{R}^{D_h \times (D-1)} \\ \mathbf{c} \in \mathbb{R}^{D_h} \\ \{\mathbf{v}_d \in \mathbb{R}^{D_h}\}_{d=2}^D \\ \{b_d \in \mathbb{R}\}_{d=1}^D \end{array} \right\} \Rightarrow D \times D_h \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \Rightarrow D \times D_h + (D-1)D_h + D$$

Number of Parameters

The total number of parameters in NADE is :

$$\left. \begin{array}{l} \mathbf{W} \in \mathbb{R}^{D_h \times (D-1)} \\ \mathbf{c} \in \mathbb{R}^{D_h} \\ \{\mathbf{v}_d \in \mathbb{R}^{D_h}\}_{d=2}^D \\ \{b_d \in \mathbb{R}\}_{d=1}^D \end{array} \right\} \Rightarrow D \times D_h \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \Rightarrow D \times D_h + (D-1)D_h + D$$

Density Estimation and Generation

Similar to FVSBN, one can generate samples from $p_\theta(\mathbf{x})$ assuming we have access to model parameters.

NADE



(a) Generated samples

NADE



(a) Generated samples



(b) NADE Possibilities ($[\lambda_1, \dots, \lambda_{784}]$)

Figure: NADE performance over BMNIST dataset (source: [4])

Generative Modeling

As usual, based on the chain rule, we have:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1)p_{\theta}(x_2|\mathbf{x}_{<2}) \dots p_{\theta}(x_d|\mathbf{x}_{) \dots p_{\theta}(x_D|\mathbf{x}_{)})$$

Generative Modeling

As usual, based on the chain rule, we have:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1)p_{\theta}(x_2|\mathbf{x}_{<2}) \dots p_{\theta}(x_d|\mathbf{x}_{) \dots p_{\theta}(x_D|\mathbf{x}_{)})$$

Now assume all the dimensions are real-valued numbers, then the modeling is:

Real NADE

Generative Modeling

As usual, based on the chain rule, we have:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1)p_{\theta}(x_2|\mathbf{x}_{<2}) \dots p_{\theta}(x_d|\mathbf{x}_{) \dots p_{\theta}(x_D|\mathbf{x}_{)})$$

Now assume all the dimensions are real-valued numbers, then the modeling is:

$$p_{\theta}(x_1) = \mathcal{N}(x_1|\mu_1, \exp(s_1))$$

Generative Modeling

As usual, based on the chain rule, we have:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1)p_{\theta}(x_2|\mathbf{x}_{<2}) \dots p_{\theta}(x_d|\mathbf{x}_{) \dots p_{\theta}(x_D|\mathbf{x}_{)})$$

Now assume all the dimensions are real-valued numbers, then the modeling is:

$$p_{\theta}(x_1) = \mathcal{N}(x_1|\mu_1, \exp(s_1))$$

⋮

$$p_{\theta}(x_d|\mathbf{x}_{)} = \mathcal{N}(x_d|\mathbf{v}_d^T \mathbf{h}_d + b_d, \exp(\mathbf{u}_d^T \mathbf{h}_d + d_d)))$$

$$\begin{cases} \mathbf{W} \in \mathbb{R}^{D_h \times (D-1)} \\ \mathbf{c}, \mathbf{v}_d, \mathbf{u}_d \in \mathbb{R}^{D_h} \\ b_d, d_d \in \mathbb{R} \end{cases}$$

Real NADE

Generative Modeling

As usual, based on the chain rule, we have:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_1)p_{\theta}(x_2|\mathbf{x}_{<2}) \dots p_{\theta}(x_d|\mathbf{x}_{) \dots p_{\theta}(x_D|\mathbf{x}_{)}$$

Now assume all the dimensions are real-valued numbers, then the modeling is:

$$p_{\theta}(x_1) = \mathcal{N}(x_1|\mu_1, \exp(s_1))$$

⋮

$$p_{\theta}(x_d|\mathbf{x}_{) = \mathcal{N}(x_d|\mathbf{v}_d^T \mathbf{h}_d + b_d, \exp(\mathbf{u}_d^T \mathbf{h}_d + d_d))$$

$$\begin{cases} \mathbf{W} \in \mathbb{R}^{D_h \times (D-1)} \\ \mathbf{c}, \mathbf{v}_d, \mathbf{u}_d \in \mathbb{R}^{D_h} \\ b_d, d_d \in \mathbb{R} \end{cases}$$

For more interesting continuous modeling, you can see [4].

Section 3

Masked Autoencoder for Distribution Estimation

Intuition From Autoencoder

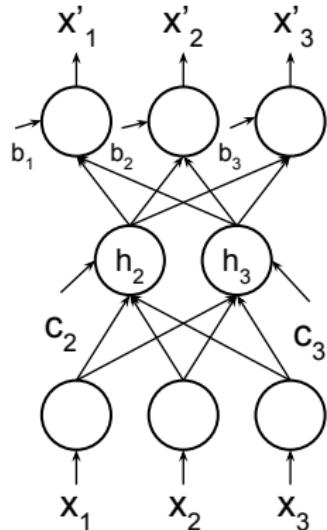
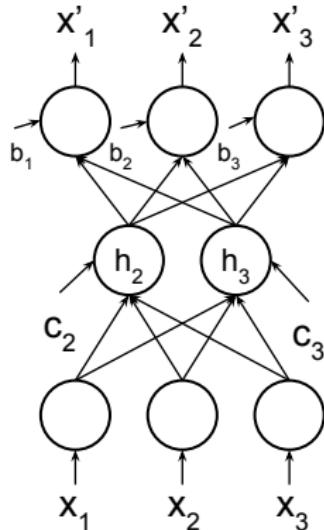


Figure: Autoencoder as unsupervised learning architecture

Intuition From Autoencoder



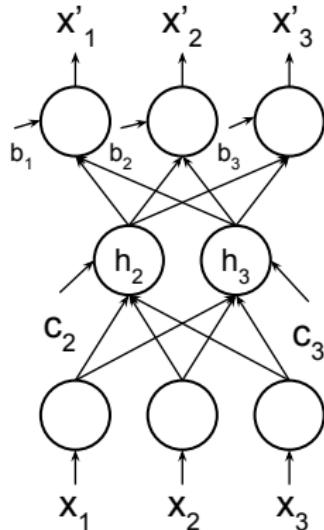
Autoencoder

An Autoencoder (AE) is characterized by:

- Encoding: $\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{c}$

Figure: Autoencoder as unsupervised learning architecture

Intuition From Autoencoder



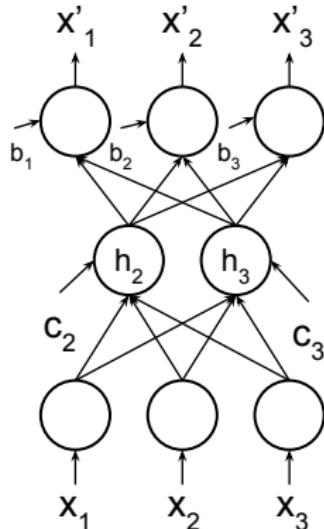
Autoencoder

An Autoencoder (AE) is characterized by:

- Encoding: $\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{c}$
- Decoding: $\mathbf{x}' = \mathbf{V}\mathbf{h} + \mathbf{b}$

Figure: Autoencoder as unsupervised learning architecture

Intuition From Autoencoder



Autoencoder

An Autoencoder (AE) is characterized by:

- Encoding: $\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{c}$
- Decoding: $\mathbf{x}' = \mathbf{V}\mathbf{h} + \mathbf{b}$
- Constraints to avoid identity mapping

Figure: Autoencoder as unsupervised learning architecture

Intuition From Autoencoder

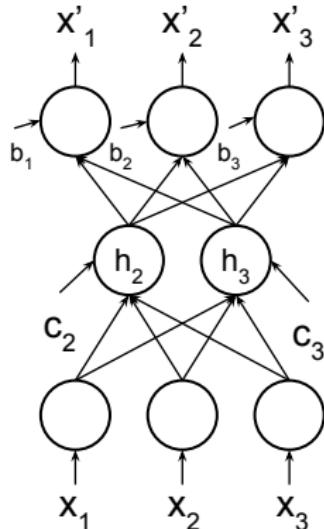


Figure: Autoencoder as unsupervised learning architecture

Autoencoder

An Autoencoder (AE) is characterized by:

- Encoding: $\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{c}$
- Decoding: $\mathbf{x}' = \mathbf{V}\mathbf{h} + \mathbf{b}$
- Constraints to avoid identity mapping
- Loss:

Binary input:

$$L_b(\boldsymbol{\theta}) = - \sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^n [x_i \log x'_i + (1 - x_i) \log(1 - x'_i)]$$

Intuition From Autoencoder

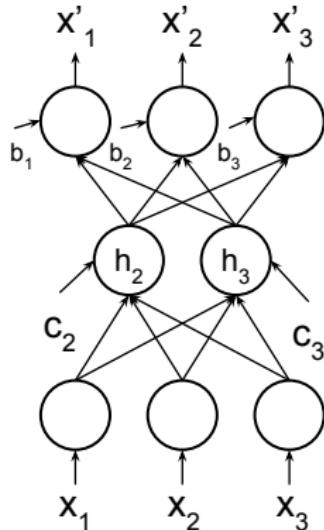


Figure: Autoencoder as unsupervised learning architecture

Autoencoder

An Autoencoder (AE) is characterized by:

- Encoding: $\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{c}$
- Decoding: $\mathbf{x}' = \mathbf{V}\mathbf{h} + \mathbf{b}$
- Constraints to avoid identity mapping
- Loss:

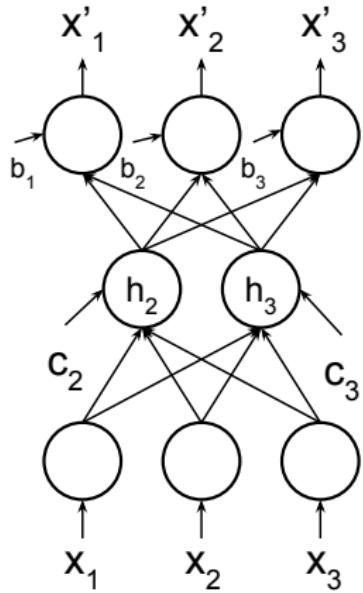
Binary input:

$$L_b(\boldsymbol{\theta}) = - \sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^n [x_i \log x'_i + (1 - x_i) \log(1 - x'_i)]$$

Continuous input:

$$L_c(\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \mathcal{D}} \sum_{i=1}^n (x_i - x'_i)^2$$

Autoencoder vs NADE



(a) Autoencoder

Autoencoder vs NADE

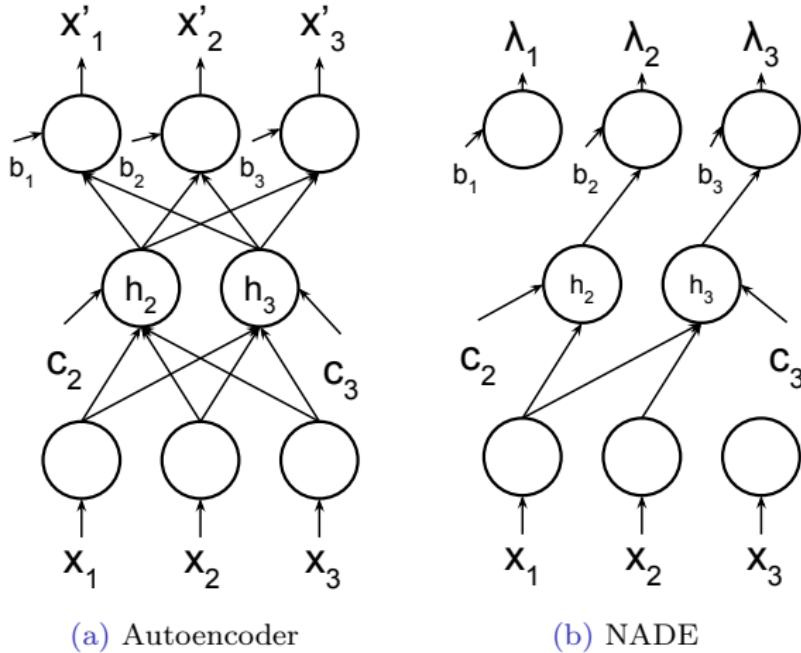


Figure: Comparison of AE and NADE (Pay attention to the architecture pruning to validate autoregressive condition)

MADE

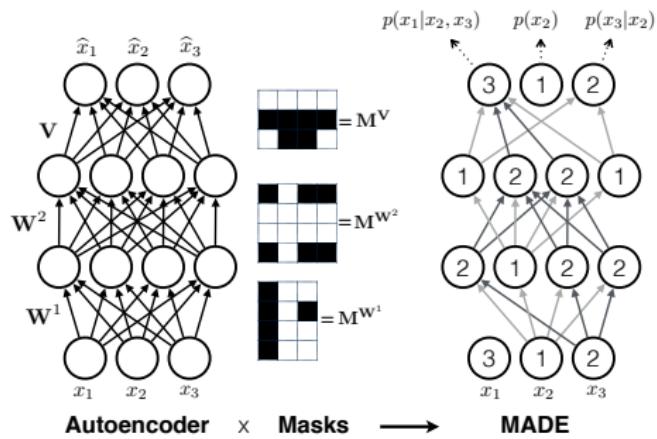


Figure: Converting an AE into a MADE
(source: [5])

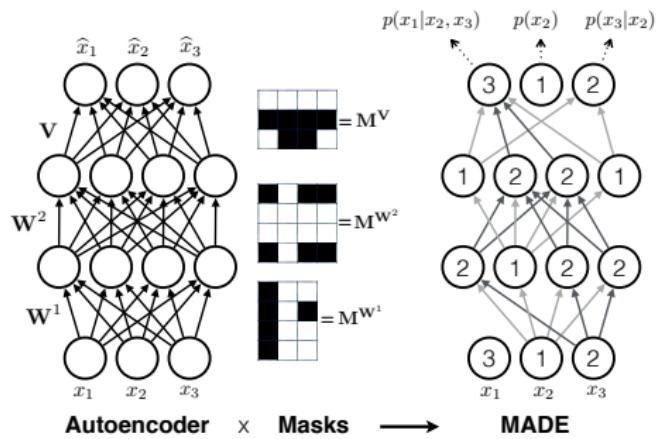
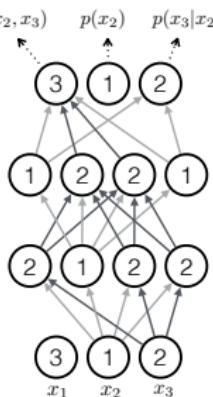


Figure: Converting an AE into a MADE
(source: [5])

AE to MADE

④ Input ordering

x_2, x_3, x_1



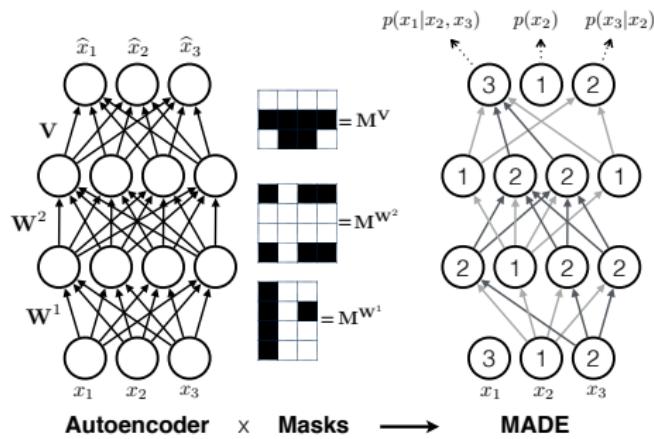


Figure: Converting an AE into a MADE
(source: [5])

AE to MADE

① Input ordering

x_2, x_3, x_1

② Set output ordering as input

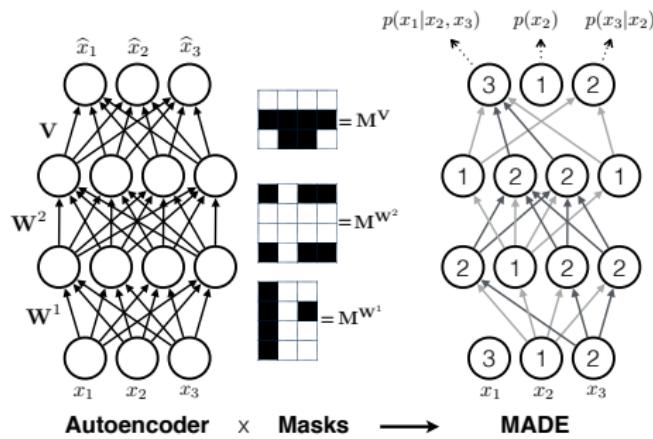


Figure: Converting an AE into a MADE
(source: [5])

AE to MADE

① Input ordering

x_2, x_3, x_1

- ② Set output ordering as input
- ③ For each unit in the hidden layer, pick a random integer index from $[1, D - 1]$

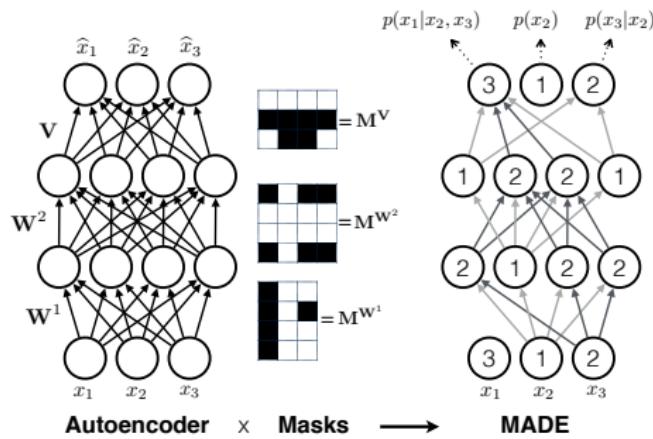


Figure: Converting an AE into a MADE
(source: [5])

AE to MADE

① Input ordering

x_2, x_3, x_1

- ② Set output ordering as input
- ③ For each unit in the hidden layer, pick a random integer index from $[1, D - 1]$
- ④ Each unit can be connected to previous layer units with equal or smaller index in hidden layers and strictly smaller index in the output layer (masks)

MADE

MADE

- In MADE you can compute all the conditional needed to evaluate $p(\mathbf{x})$ in parallel.

MADE

- In MADE you can compute all the conditional needed to evaluate $p(\mathbf{x})$ in parallel.
- Sampling from MADE is sequentially

MADE

- In MADE you can compute all the conditional needed to evaluate $p(\mathbf{x})$ in parallel.
- Sampling from MADE is sequentially
- Assume ordering x_2, x_3, x_1 , then the chain rule is:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_2)p_{\theta}(x_3|x_2)p_{\theta}(x_1|x_2, x_3)$$

MADE

- In MADE you can compute all the conditionals needed to evaluate $p(\mathbf{x})$ in parallel.
- Sampling from MADE is sequentially
- Assume ordering x_2, x_3, x_1 , then the chain rule is:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_2)p_{\theta}(x_3|x_2)p_{\theta}(x_1|x_2, x_3)$$

- Because distribution over x_2 is unconditional, the second unit in the output layer has no input connection.

MADE

- In MADE you can compute all the conditionals needed to evaluate $p(\mathbf{x})$ in parallel.
- Sampling from MADE is sequentially
- Assume ordering x_2, x_3, x_1 , then the chain rule is:

$$p_{\theta}(\mathbf{x}) = p_{\theta}(x_2)p_{\theta}(x_3|x_2)p_{\theta}(x_1|x_2, x_3)$$

- Because distribution over x_2 is unconditional, the second unit in the output layer has no input connection.
- Because x_1 is in none of the conditions, the first unit in the input layer has no output connection.

Section 4

Recurrent Neural Networks

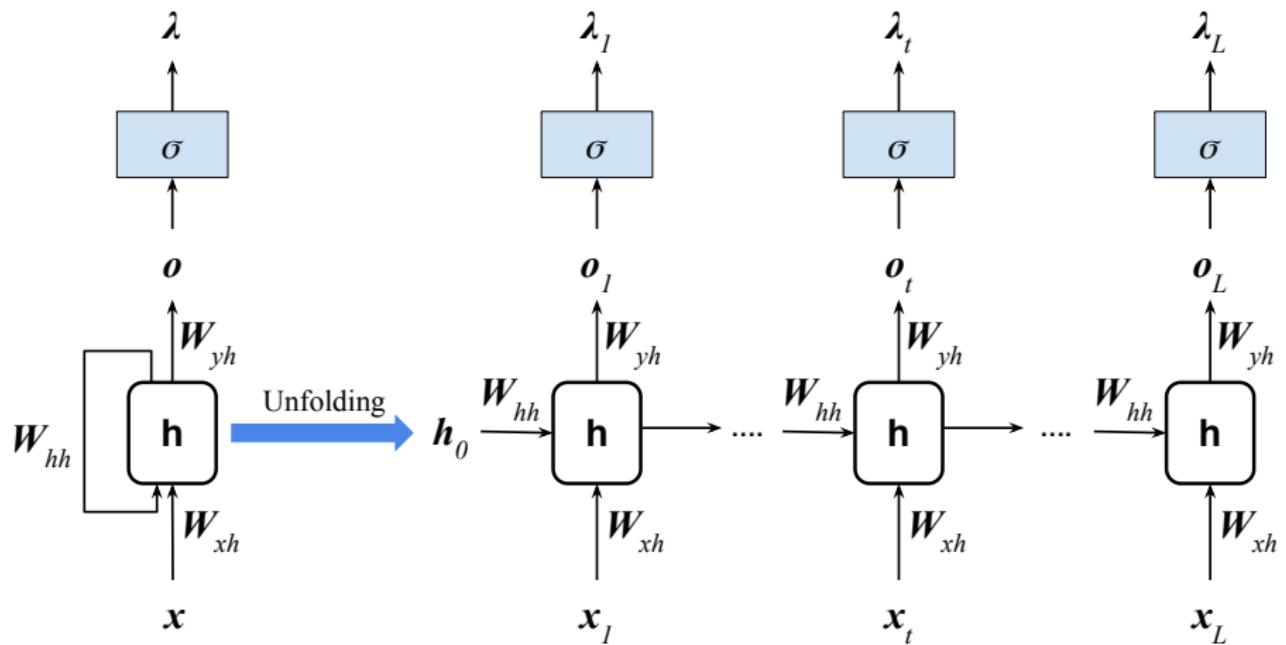


Figure: RNN architecture with its unfolded version

RNN

Assume: $x = \overbrace{\text{I}}^{\wedge} \quad \overbrace{\text{am}}^{\wedge} \quad \overbrace{\text{going}}^{\wedge} \quad \overbrace{\text{to}}^{\wedge} \quad \overbrace{\text{Bank}}^{\wedge}$

RNN

Assume: $x = \overbrace{\text{I}}^{x_1, y_1=0} \quad \overbrace{\text{am}} \quad \overbrace{\text{going}} \quad \overbrace{\text{to}} \quad \overbrace{\text{Bank}}$

RNN

Assume: $x = \overbrace{\text{I}}^{x_1, y_1=0} \overbrace{\text{am}}^{x_2, y_2=0} \overbrace{\text{going}} \overbrace{\text{to}} \overbrace{\text{Bank}}$

RNN

Assume: $x = \overbrace{\text{I}}^{x_1, y_1=0} \overbrace{\text{am}}^{x_2, y_2=0} \overbrace{\text{going}}^{x_3, y_3=1} \overbrace{\text{to}} \quad \overbrace{\text{Bank}}$

RNN

Assume: $x = \overbrace{\text{I}}^{x_1, y_1=0} \overbrace{\text{am}}^{x_2, y_2=0} \overbrace{\text{going}}^{x_3, y_3=1} \overbrace{\text{to}}^{x_4, y_4=0} \overbrace{\text{Bank}}$

RNN

Assume: $x = \overbrace{\text{I}}^{x_1, y_1=0} \overbrace{\text{am}}^{x_2, y_2=0} \overbrace{\text{going}}^{x_3, y_3=1} \overbrace{\text{to}}^{x_4, y_4=0} \overbrace{\text{Bank}}^{x_3, y_3=0}$

RNN

Assume: $\mathbf{x} = \overbrace{\text{I}}^{x_1, y_1=0} \overbrace{\text{am}}^{x_2, y_2=0} \overbrace{\text{going}}^{x_3, y_3=1} \overbrace{\text{to}}^{x_4, y_4=0} \overbrace{\text{Bank}}^{x_3, y_3=0}$

Thus you have problem to relate two sequences of length l where l is variable and your dataset is $\mathcal{D} = \{(\mathbf{X}_i \in \mathbb{R}^{D \times l_i}, \mathbf{y}_i \in \mathbb{R}^{l_i})\}_{i=1}^N$. RNNs can capture this scenario.

RNN

- Initialization: $\mathbf{h}_0 = \mathbf{0}$
- For $i = 1$ to l :
 - Hidden state update: $\mathbf{h}_i = f(\mathbf{W}_{hh}\mathbf{h}_{i-1} + \mathbf{W}_{xh}\mathbf{x}_i)$
 - Output calculation: $\mathbf{o}_i = \sigma(\mathbf{W}_{yh}\mathbf{h}_i)$

Generative RNN

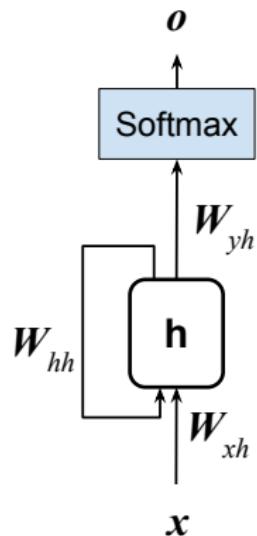


Figure: Compact generative RNN

Generative RNN

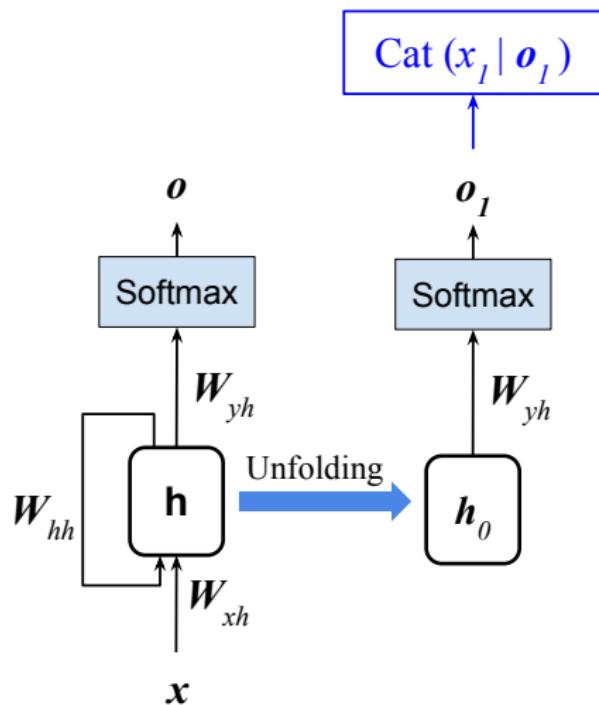


Figure: Calculations for $p_\theta(x_1)$

Generative RNN

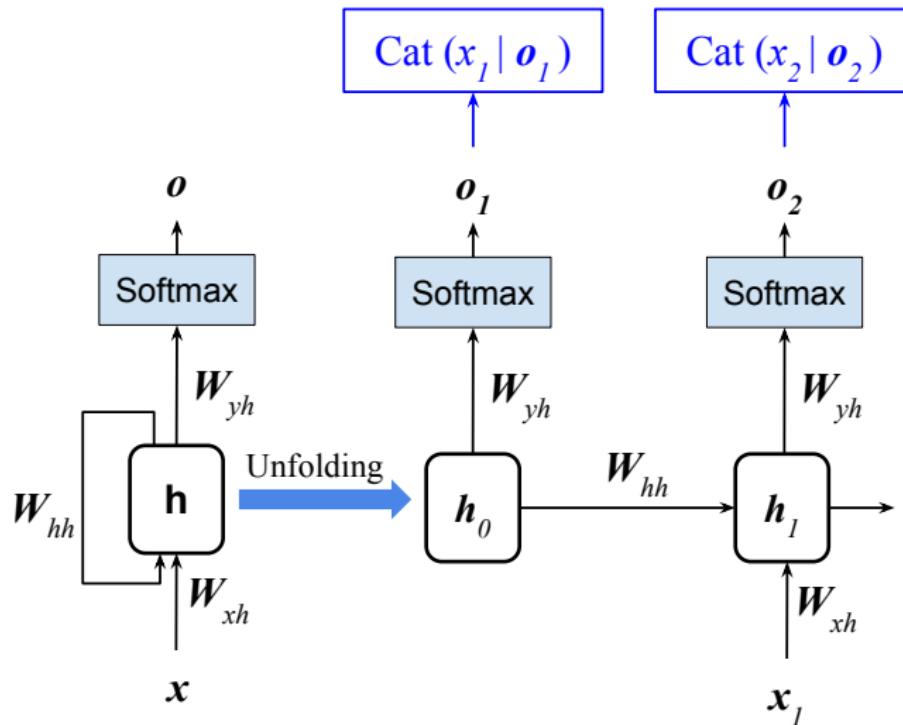


Figure: Calculations for $p_\theta(x_2|x_1)$

Generative RNN

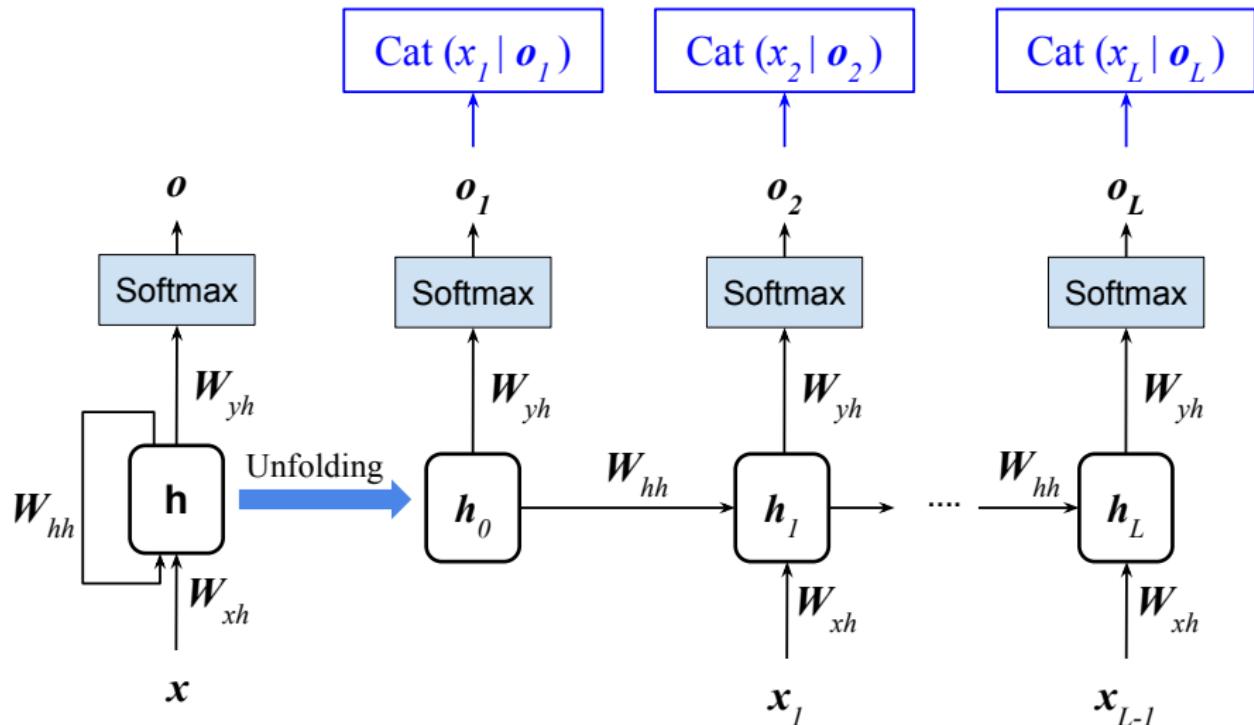


Figure: Calculations for $p_\theta(x_D | x_1, \dots, x_{D-1})$

Generative Modeling Based on RNN Architecture

Pixel RNN

Assume we want to design an autoregressive generative model for MNIST dataset $\mathcal{D} = \{\mathbf{x}_i \in \{0, 1, \dots, 255\}^{784}\}_{i=1}^N$ using RNN architecture:

$$\mathbf{h}_0 = \mathbf{b}_0 \Rightarrow \mathbf{o}_0 = \text{softmax}(\mathbf{W}_{yh}\mathbf{h}_0) \Rightarrow p(x_1) = \text{Cat}(x_1|\mathbf{o}_1)$$

Generative Modeling Based on RNN Architecture

Pixel RNN

Assume we want to design an autoregressive generative model for MNIST dataset $\mathcal{D} = \{\mathbf{x}_i \in \{0, 1, \dots, 255\}^{784}\}_{i=1}^N$ using RNN architecture:

$$\mathbf{h}_0 = \mathbf{b}_0 \Rightarrow \mathbf{o}_0 = \text{softmax}(\mathbf{W}_{yh}\mathbf{h}_0) \Rightarrow p(x_1) = \text{Cat}(x_1|\mathbf{o}_1)$$

⋮

$$\mathbf{h}_{i-1} = f(\mathbf{W}_{hh}\mathbf{h}_{i-2} + \mathbf{W}_{xh}\mathbf{x}_{i-1}) \Rightarrow$$

$$\mathbf{o}_{i-1} = \text{softmax}(\mathbf{W}_{yh}\mathbf{h}_{i-1}) \Rightarrow p(x_i|\mathbf{x}_{<i}) = \text{Cat}(x_i|\mathbf{o}_{i-1})$$

Generative Modeling Based on RNN Architecture

Pixel RNN

Assume we want to design an autoregressive generative model for MNIST dataset $\mathcal{D} = \{\mathbf{x}_i \in \{0, 1, \dots, 255\}^{784}\}_{i=1}^N$ using RNN architecture:

$$\mathbf{h}_0 = \mathbf{b}_0 \Rightarrow \mathbf{o}_0 = \text{softmax}(\mathbf{W}_{yh}\mathbf{h}_0) \Rightarrow p(x_1) = \text{Cat}(x_1|\mathbf{o}_1)$$

⋮

$$\mathbf{h}_{i-1} = f(\mathbf{W}_{hh}\mathbf{h}_{i-2} + \mathbf{W}_{xh}\mathbf{x}_{i-1}) \Rightarrow$$

$$\mathbf{o}_{i-1} = \text{softmax}(\mathbf{W}_{yh}\mathbf{h}_{i-1}) \Rightarrow p(x_i|\mathbf{x}_{<i}) = \text{Cat}(x_i|\mathbf{o}_{i-1})$$

Thus we can calculate conditionals while RNN guarantees that the autoregression constraint is met.

Generative RNN

Generative RNN

- Pros:

- Variable-length input
- Fixed model size
- History abstraction via hidden representation \mathbf{h}_d
- Weight sharing

- Cons

- Sequential computations
- Difficult to carry information for long sequences
- Ordering required

Pixel RNN



Figure: Pixel RNN results in image completion (source: [6])

Section 5

Learning

Distance Metric

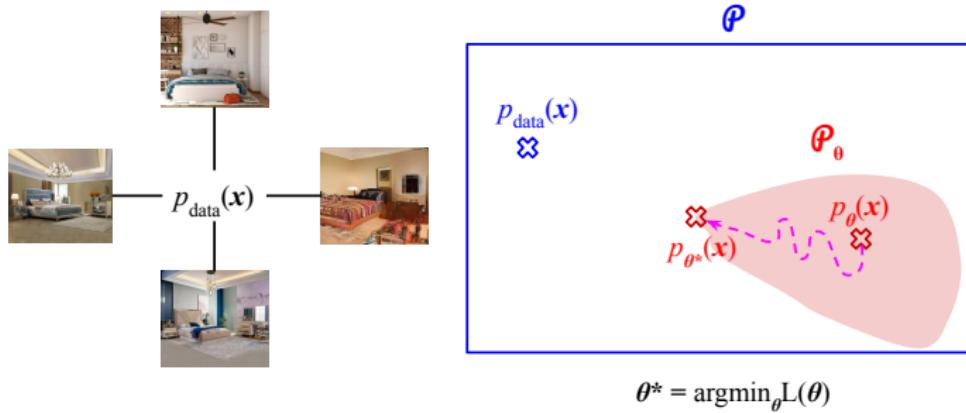


Figure: General view to generative modeling

Distance Metric

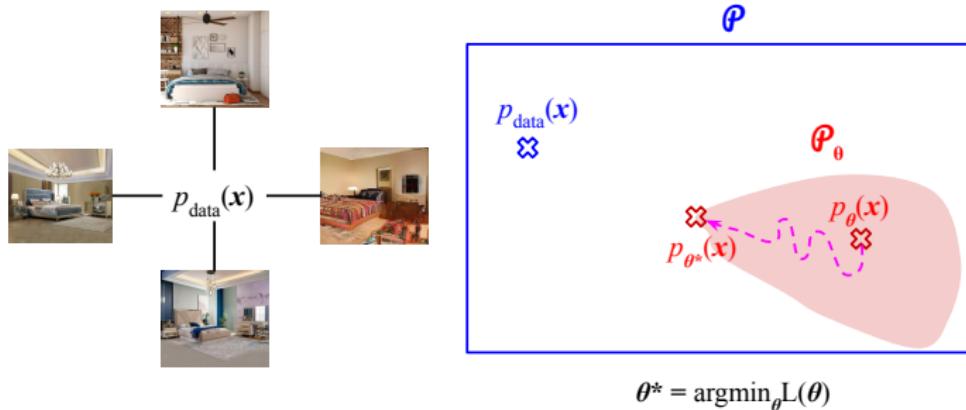


Figure: General view to generative modeling

So Far

- Characterizing \mathcal{P}
- Characterizing \mathcal{P}_θ
 - FVSBN, NADE, MADE, RNN, ...

Distance Metric

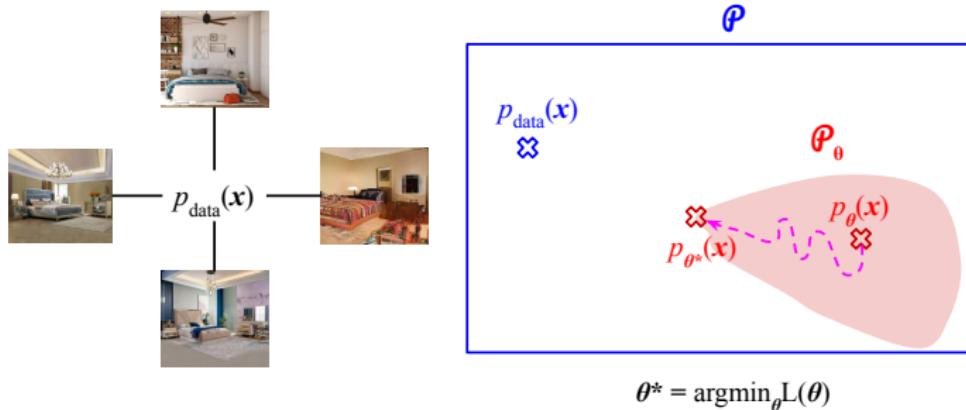


Figure: General view to generative modeling

So Far

- Characterizing \mathcal{P}
- Characterizing \mathcal{P}_θ
 - FVSBN, NADE, MADE, RNN, ...

From Now on

- Determining the distance $L(\theta)$
- Optimizing θ (Finding θ^*)

Distance Metric

Distance Metric

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

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

Distance Metric

Distance Metric

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

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

Distance Metric

Distance Metric

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

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

Distance Metric

Distance Metric

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

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

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

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

Distance Metric

Distance Metric

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

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

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

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

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

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

From KL divergence to Model Likelihood

Model Likelihood

We see:

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

Thus:

From KL divergence to Model Likelihood

Model Likelihood

We see:

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

Thus:

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

From KL divergence to Model Likelihood

Model Likelihood

We see:

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

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_{data}

From KL divergence to Model Likelihood

Model Likelihood

We see:

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

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_{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:

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

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

Model Likelihood Estimation

Model Likelihood Estimation

We are interested in solving the following problem:

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

but we don't have access to p_{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:

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

Thus:

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

Optimization

Optimization

By adding a minus sign to the objective function, we have the following minimization problem:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} -\frac{1}{|\mathcal{D}|} \sum_{n=1}^N \log p_{\boldsymbol{\theta}}(\mathbf{x}_n)$$

and we can use different methods to find $\boldsymbol{\theta}^*$ such as:

- Gradient descent
- Stochastic gradient descent
- ...

Section 6

Examples

Coin Tossing

Problem

Assume you have a biased coin and toss it for 4 times and your dataset is:

$$\mathcal{D} = \{1, 0, 1, 1\}, \begin{cases} 1 & \text{head} \\ 0 & \text{tail} \end{cases}$$

Suggest the parameterized model and find the optimum parameters.

Coin Tossing

Problem

Assume you have a biased coin and toss it for 4 times and your dataset is:

$$\mathcal{D} = \{1, 0, 1, 1\}, \begin{cases} 1 & \text{head} \\ 0 & \text{tail} \end{cases}$$

Suggest the parameterized model and find the optimum parameters.

Solution

The model is a Bernoulli model as: $p_\theta(x) = \theta^x(1 - \theta)^{1-x}$

The optimum parameter is the result of the following optimization problem:

$$\theta^* = \operatorname{argmin}_\theta -\frac{1}{4} \sum_{i=1}^4 [x_i \log \theta + (1 - x_i) \log(1 - \theta)]$$

Coin Tossing

Solution (Cont.)

Taking into account the dataset $\mathcal{D} = \{1, 0, 1, 1\}$, we have:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} -\frac{1}{4}(\log \theta + \log(1-\theta) + \log \theta + \log \theta) \\ &= -\frac{3}{4} \log \theta - \frac{1}{4} \log(1-\theta)\end{aligned}$$

Coin Tossing

Solution (Cont.)

Taking into account the dataset $\mathcal{D} = \{1, 0, 1, 1\}$, we have:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} -\frac{1}{4}(\log \theta + \log(1-\theta) + \log \theta + \log \theta) \\ &= -\frac{3}{4} \log \theta - \frac{1}{4} \log(1-\theta)\end{aligned}$$

By taking the derivative of the objective with respect to θ , we can find:

$$\theta^* = \frac{3}{4}$$

Coin Tossing

Solution (Cont.)

Taking into account the dataset $\mathcal{D} = \{1, 0, 1, 1\}$, we have:

$$\begin{aligned}\theta^* &= \operatorname{argmin}_{\theta} -\frac{1}{4}(\log \theta + \log(1-\theta) + \log \theta + \log \theta) \\ &= -\frac{3}{4} \log \theta - \frac{1}{4} \log(1-\theta)\end{aligned}$$

By taking the derivative of the objective with respect to θ , we can find:

$$\theta^* = \frac{3}{4}$$

Thus your generative model is: $p_{\theta}(x) = \text{Ber}(x|\frac{3}{4})$ and you can:

- Sample from this distribution
- Calculate the probability of an input x

Coin Tossing

General Case

Now assume you have access to a dataset with N samples where N_h samples are head and N_t samples are tail ($N = N_h + N_t$), then we can write:

$$\theta^* = \operatorname{argmin}_{\theta} - \left(\frac{N_h}{N} \log \theta + \frac{N_t}{N} \log(1 - \theta) \right)$$

Coin Tossing

General Case

Now assume you have access to a dataset with N samples where N_h samples are head and N_t samples are tail ($N = N_h + N_t$), then we can write:

$$\theta^* = \operatorname{argmin}_{\theta} - \left(\frac{N_h}{N} \log \theta + \frac{N_t}{N} \log(1 - \theta) \right)$$

which results in:

$$\theta^* = \frac{N_h}{N}$$

Coin Tossing

General Case

Now assume you have access to a dataset with N samples where N_h samples are head and N_t samples are tail ($N = N_h + N_t$), then we can write:

$$\theta^* = \operatorname{argmin}_{\theta} - \left(\frac{N_h}{N} \log \theta + \frac{N_t}{N} \log(1 - \theta) \right)$$

which results in:

$$\theta^* = \frac{N_h}{N}$$

and thus your generative model is:

$$p_{\theta}(x) = \text{Ber} \left(x \mid \frac{N_h}{N} \right)$$

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$.

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. Then:

- Select a (mini)batch \mathcal{B} of samples with size N' .

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. Then:

- Select a (mini)batch \mathcal{B} of samples with size N' .
- Compute the following loss function in a parametric form:

$$L(\boldsymbol{\theta}) = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. Then:

- Select a (mini)batch \mathcal{B} of samples with size N' .
- Compute the following loss function in a parametric form:

$$L(\boldsymbol{\theta}) = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

- Take the gradient with respect to parameters $\nabla_{\boldsymbol{\theta}} L$

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. Then:

- Select a (mini)batch \mathcal{B} of samples with size N' .
- Compute the following loss function in a parametric form:

$$L(\boldsymbol{\theta}) = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

- Take the gradient with respect to parameters $\nabla_{\boldsymbol{\theta}} L$
- Update parameters using (stochastic) gradient descent or any other gradient-based optimization method

Learning General Autoregressive Model

Learning

Assume you're given an Autoregressive model (such as NADE, MADE, Pixel-RNN, etc.), and you are asked to train the model using $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$. Then:

- Select a (mini)batch \mathcal{B} of samples with size N' .
- Compute the following loss function in a parametric form:

$$L(\boldsymbol{\theta}) = -\frac{1}{|\mathcal{B}|} \sum_{\mathbf{x}_i \in \mathcal{B}} \log p_{\boldsymbol{\theta}}(\mathbf{x}_i)$$

- Take the gradient with respect to parameters $\nabla_{\boldsymbol{\theta}} L$
- Update parameters using (stochastic) gradient descent or any other gradient-based optimization method
- Repeat the above steps until a convergence criterion is met.

Avoid Overfitting

Avoid Overfitting

The following strategies can be used to avoid overfitting:

Avoid Overfitting

Avoid Overfitting

The following strategies can be used to avoid overfitting:

- Decreasing the size of $p_\theta(\mathbb{X})$ set (Toward independence assumption) by:
 - Limiting the number of processors that each variable can depend on:

$$p_\theta(x_i | \mathbf{x}_{<i}) = p_\theta(x_i | x_{i-1}, x_{i-2})$$

Avoid Overfitting

Avoid Overfitting

The following strategies can be used to avoid overfitting:

- Decreasing the size of $p_{\theta}(\mathbb{X})$ set (Toward independence assumption) by:
 - Limiting the number of processors that each variable can depend on:

$$p_{\theta}(x_i | \mathbf{x}_{<i}) = p_{\theta}(x_i | x_{i-1}, x_{i-2})$$

- Weight sharing (similar to NADE and RNN)

Avoid Overfitting

Avoid Overfitting

The following strategies can be used to avoid overfitting:

- Decreasing the size of $p_{\theta}(\mathbb{X})$ set (Toward independence assumption) by:
 - Limiting the number of processors that each variable can depend on:

$$p_{\theta}(x_i | \mathbf{x}_{<i}) = p_{\theta}(x_i | x_{i-1}, x_{i-2})$$

- Weight sharing (similar to NADE and RNN)
- Regularization such as Gaussian prior over parameters (Weight decay)

Avoid Overfitting

Avoid Overfitting

The following strategies can be used to avoid overfitting:

- Decreasing the size of $p_{\theta}(\mathbb{X})$ set (Toward independence assumption) by:
 - Limiting the number of processors that each variable can depend on:

$$p_{\theta}(x_i | \mathbf{x}_{<i}) = p_{\theta}(x_i | x_{i-1}, x_{i-2})$$

- Weight sharing (similar to NADE and RNN)
- Regularization such as Gaussian prior over parameters (Weight decay)
- Early stopping via monitoring likelihood over a validation set

Section 7

Conditional Generation

Conditional Autoregressive Generation

Conditional Generation Vs. Supervised Learning

Assume dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ where \mathbf{y}_i is a text and \mathbf{x}_i is the corresponding speech. For this problem:

Conditional Autoregressive Generation

Conditional Generation Vs. Supervised Learning

Assume dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ where \mathbf{y}_i is a text and \mathbf{x}_i is the corresponding speech. For this problem:

- Similar to supervised learning, we are interested in modeling $p(\mathbb{X}|\mathbf{y})$. Thus the unconditional distribution $p(\mathbb{Y})$ is not modeled and \mathbf{y} is assumed to be given.

Conditional Autoregressive Generation

Conditional Generation Vs. Supervised Learning

Assume dataset $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ where \mathbf{y}_i is a text and \mathbf{x}_i is the corresponding speech. For this problem:

- Similar to supervised learning, we are interested in modeling $p(\mathbb{X}|\mathbf{y})$. Thus the unconditional distribution $p(\mathbb{Y})$ is not modeled and \mathbf{y} is assumed to be given.
- In contrast to supervised learning, \mathbb{X} is a high-dimensional random vector (note that for binary classification X is Bernoulli and for multi-class classification, X is Categorical).

Conditional Autoregressive Generation

Modeling

Keeping in mind that $\mathbf{y} \in \mathbb{R}^{D_g}$ is a give vector and we want to model $p_\theta(\mathbb{X}|\mathbf{y})$ using Autoregressive model, where \mathbb{X} is D dimensional random vector, we have:

Conditional Autoregressive Generation

Modeling

Keeping in mind that $\mathbf{y} \in \mathbb{R}^{D_g}$ is a give vector and we want to model $p_\theta(\mathbb{X}|\mathbf{y})$ using Autoregressive model, where \mathbb{X} is D dimensional random vector, we have:

$$p_\theta(\mathbb{X}|\mathbf{y}) = p_\theta(X_1|\mathbf{y})p_\theta(X_2|\mathbf{x}_{<2}, \mathbf{y}) \dots p_\theta(X_d|\mathbf{x}_{, \mathbf{y}}) \dots p_\theta(X_D|\mathbf{x}_{, \mathbf{y}})$$

and we can use any of the autoregressive modeling approaches to model the conditionals.

Conditional Autoregressive Generation

Modeling

Keeping in mind that $\mathbf{y} \in \mathbb{R}^{D_g}$ is a give vector and we want to model $p_\theta(\mathbb{X}|\mathbf{y})$ using Autoregressive model, where \mathbb{X} is D dimensional random vector, we have:

$$p_\theta(\mathbb{X}|\mathbf{y}) = p_\theta(X_1|\mathbf{y})p_\theta(X_2|\mathbf{x}_{<2}, \mathbf{y}) \dots p_\theta(X_d|\mathbf{x}_{, \mathbf{y}}) \dots p_\theta(X_D|\mathbf{x}_{, \mathbf{y}})$$

and we can use any of the autoregressive modeling approaches to model the conditionals.

Density Estimation and Generation

Similar to unconditional Autoregressive models, you can sample new \mathbf{x} given \mathbf{y} (assume different response \mathbf{x} that can ChatGPT give to a fixed question \mathbf{y}), or estimate density for a \mathbf{x} vector while \mathbf{y} is given.

Learning

Distance Metric

One option for distance metric is:

Learning

Distance Metric

One option for distance metric is:

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

Learning

Distance Metric

One option for distance metric is:

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

Learning

Distance Metric

One option for distance metric is:

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

Learning

Distance Metric

One option for distance metric is:

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

Learning

Distance Metric

One option for distance metric is:

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

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.

Learning

Distance Metric

Altogether:

Learning

Distance Metric

Altogether:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbb{Y})} \left[\text{KL} \left(p_{\text{data}}(\mathbb{X}|\mathbf{y}) \parallel p_{\boldsymbol{\theta}}(\mathbb{X}|\mathbf{y}) \right) \right]$$

Learning

Distance Metric

Altogether:

$$\begin{aligned}\theta^* &= \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbb{Y})} \left[\text{KL} \left(p_{\text{data}}(\mathbb{X} | \mathbf{y}) \parallel p_{\theta}(\mathbb{X} | \mathbf{y}) \right) \right] \\ &\equiv \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}(\mathbb{X}, \mathbb{Y})} \left[\log p_{\theta}(\mathbf{x} | \mathbf{y}) \right]\end{aligned}$$

Learning

Distance Metric

Altogether:

$$\begin{aligned}\boldsymbol{\theta}^* &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbb{Y})} \left[\text{KL} \left(p_{\text{data}}(\mathbb{X}|\mathbf{y}) \parallel p_{\boldsymbol{\theta}}(\mathbb{X}|\mathbf{y}) \right) \right] \\ &\equiv \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p_{\text{data}}(\mathbb{X}, \mathbb{Y})} \left[\log p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{y}) \right]\end{aligned}$$

The above can be translated as:

Minimizing the expected KL divergence of true conditional distribution and model conditional distribution over true \mathbb{X} is equivalent to maximizing the expected likelihood of model conditional distribution over true joint distribution.

Model Likelihood Estimation

Model Conditional Likelihood Estimation

We are interested in solving the following problem:

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

but we don't have access to $p_{\text{data}}(\mathbb{X}, \mathbb{Y})$

Model Likelihood Estimation

Model Conditional Likelihood Estimation

We are interested in solving the following problem:

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

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

Model Likelihood Estimation

Model Conditional Likelihood Estimation

We are interested in solving the following problem:

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

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

Solution via Monte Carlo Estimate

Using the Monte Carlo estimate we have:

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

Model Likelihood Estimation

Model Conditional Likelihood Estimation

We are interested in solving the following problem:

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

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

Solution via Monte Carlo Estimate

Using the Monte Carlo estimate we have:

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

Thus:

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

Section 8

Applications

Adversarial Robustness

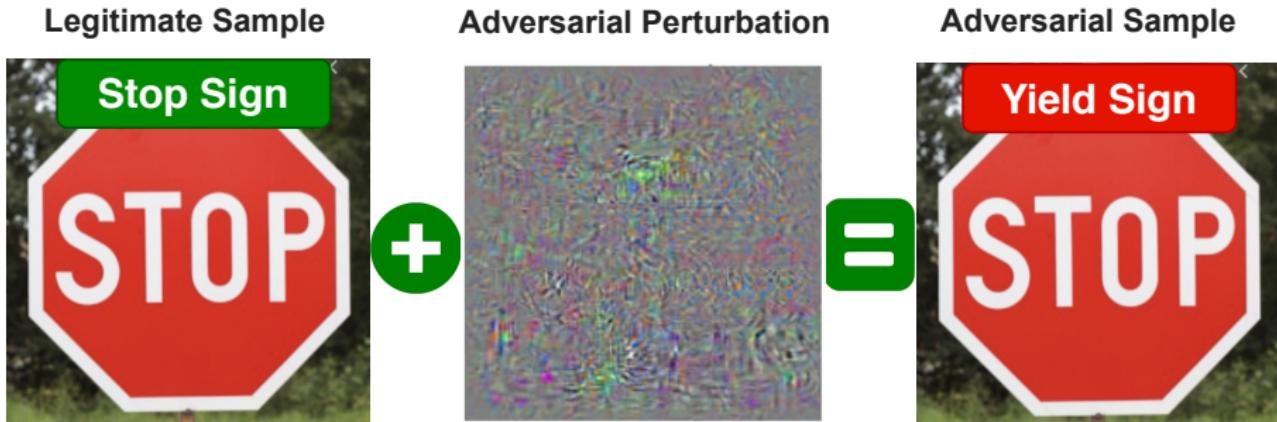


Figure: Different adversarial attacks to Frog image from Cifar10 dataset (source: [7])

Adversarial Robustness



Figure: Sample adversarial attack to deep learning architectures (source: [7])

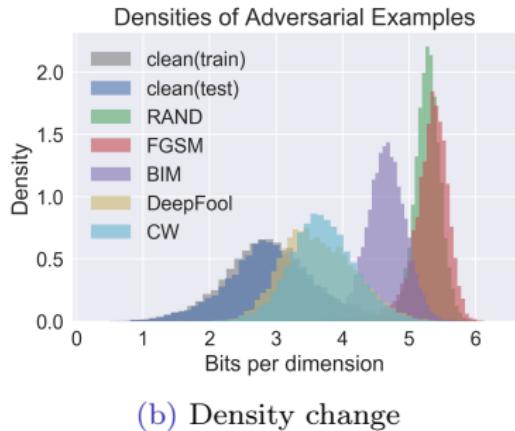
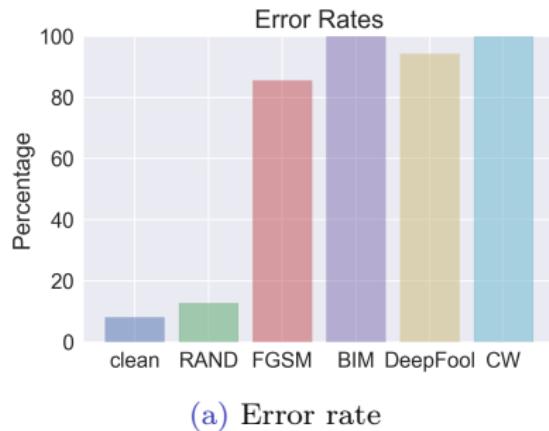


Figure: Using autoregressive models to detect adversarial samples (source: [7])

Section 9

Summary

Summary

Idea

- Using chain rule to factorize the joint

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:
 - Easy probability estimation

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:
 - Easy probability estimation
 - Easy to train and monitor training via validation likelihood
- Cons:

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:
 - Easy probability estimation
 - Easy to train and monitor training via validation likelihood
- Cons:
 - Ordering requirement

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:
 - Easy probability estimation
 - Easy to train and monitor training via validation likelihood
- Cons:
 - Ordering requirement
 - Sequential generation

Summary

Idea

- Using chain rule to factorize the joint
- Parameterizing the factors in the chain rule using neural networks
- Using KL divergence as the distance metric
- Using Monte Carlo estimation to approximate the distance metric

Pros and Cons

- Pros:
 - Easy probability estimation
 - Easy to train and monitor training via validation likelihood
- Cons:
 - Ordering requirement
 - Sequential generation
 - No high-level representation z

List of Abbreviations

Complete	Abbreviation
Binary Logistic Regression	BLR
Fully Visible Sigmoid Belief Network	FVSBN
Masked Autoencoder for Distribution Estimation	MADE
Neural Autoregressive Density Estimation	NADE
Recurrent Neural Networks	RNN

References I

-  Kevin P Murphy,
Machine learning: a probabilistic perspective,
MIT press, 2012.
-  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.
-  Kevin P Murphy,
Probabilistic machine learning: an introduction,
MIT press, 2022.
-  Benigno Uria, Marc-Alexandre Côté, Karol Gregor, Iain Murray, and Hugo Larochelle,
“Neural autoregressive distribution estimation,”
The Journal of Machine Learning Research, vol. 17, no. 1, pp. 7184–7220, 2016.
-  Mathieu Germain, Karol Gregor, Iain Murray, and Hugo Larochelle,
“Made: Masked autoencoder for distribution estimation,”
in *International conference on machine learning*. PMLR, 2015, pp. 881–889.
-  Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu,
“Pixel recurrent neural networks,”
in *International conference on machine learning*. PMLR, 2016, pp. 1747–1756.

References II



Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman,
“Pixeldefend: Leveraging generative models to understand and defend against adversarial
examples,”
arXiv preprint arXiv:1710.10766, 2017.