

Section 6

From Supervised Learning to Generative Modeling

Subsection 1

Logistic Regression

Logistic Regression Model

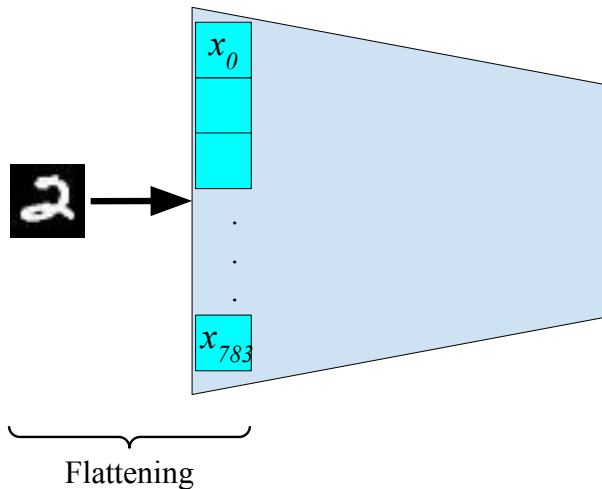


Figure: Logistic regression steps

Logistic Regression Model

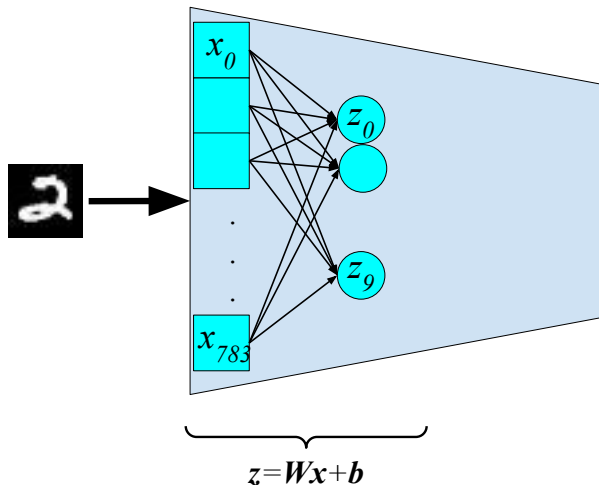


Figure: Logistic regression steps

Logistic Regression Model

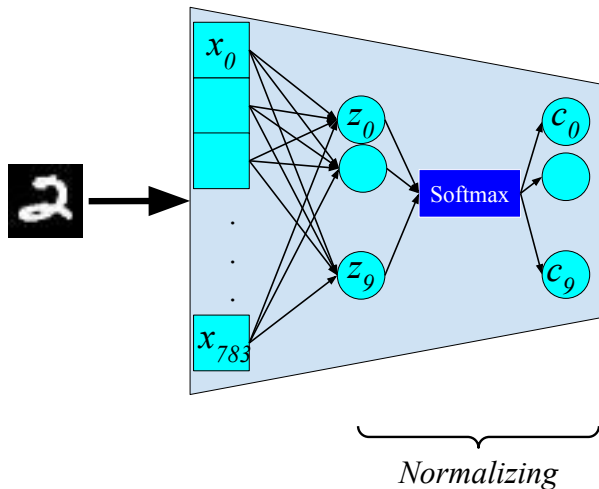
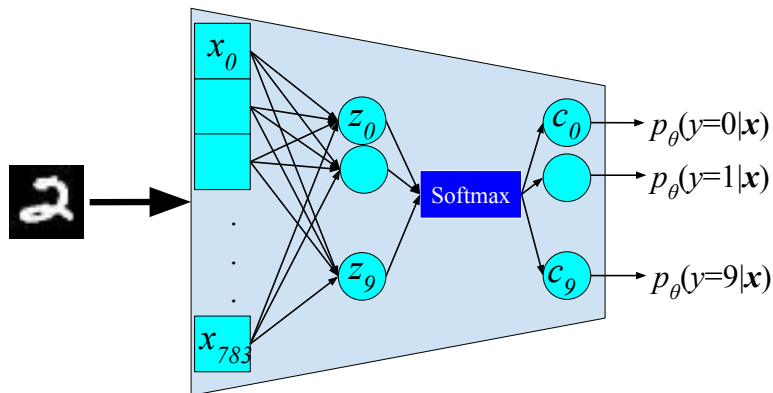


Figure: Logistic regression steps

Logistic Regression Model



$$\underbrace{\hspace{10em}}_{p_{\theta}(y|\mathbf{x}) = \text{Cat}(y;\mathbf{C})}$$

Figure: Logistic regression steps

Logistic Regression Model

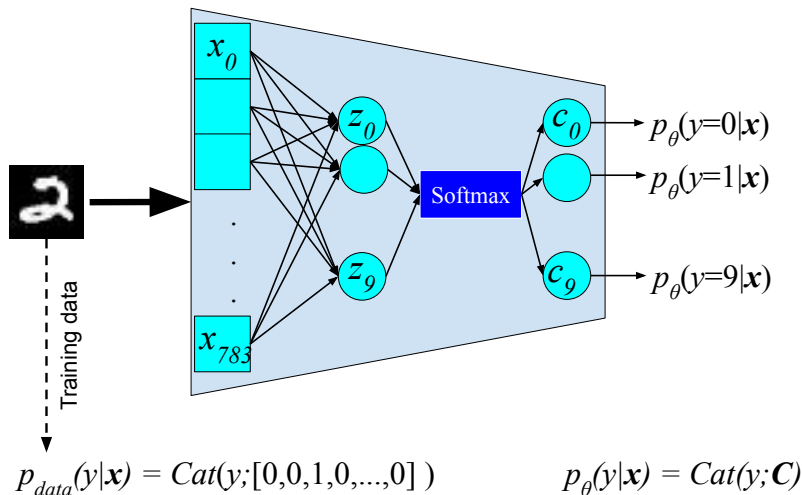


Figure: Logistic regression steps

Logistic Regression Model

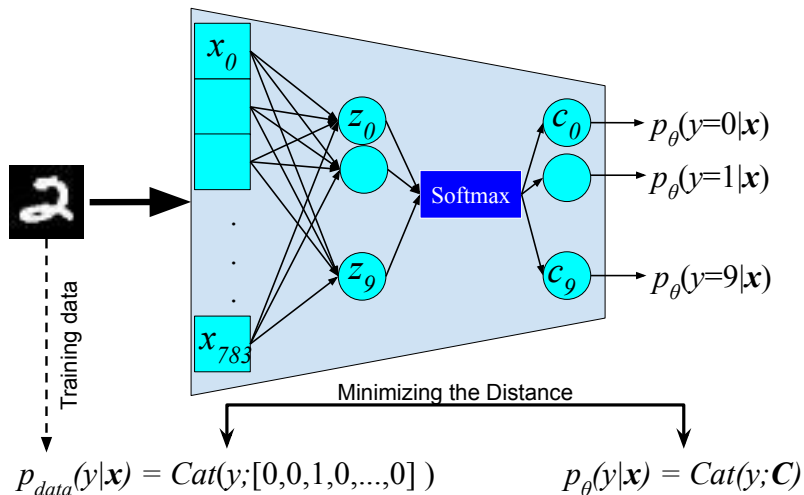


Figure: Logistic regression steps

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:

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

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:

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

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:

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

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

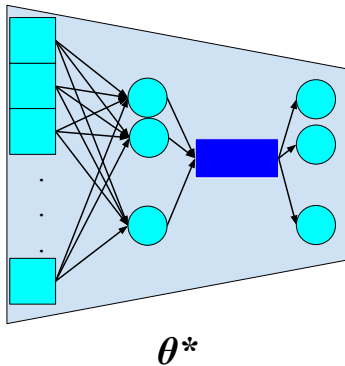


Figure: Sampling a trained model

Sampling

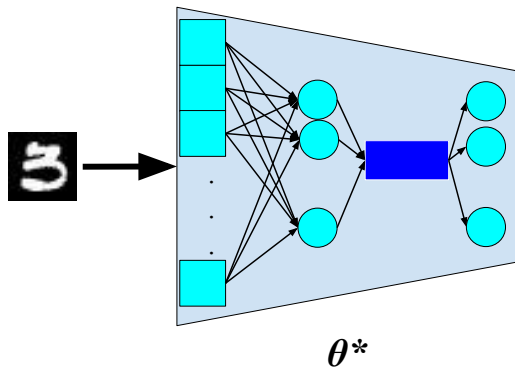


Figure: Sampling a trained model

Sampling

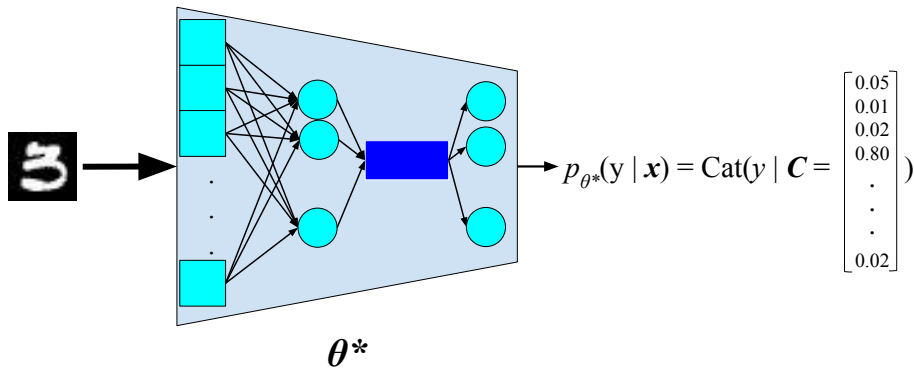


Figure: Sampling a trained model

Sampling

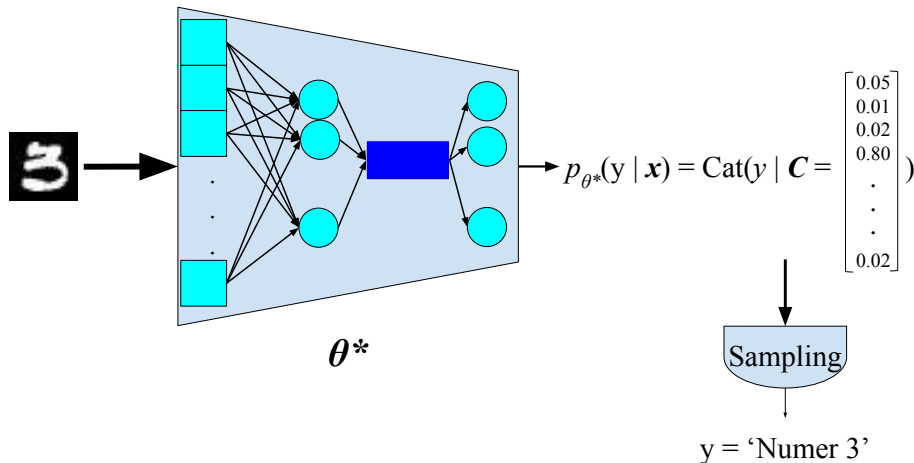


Figure: Sampling a trained model

Sampling a Categorical Distribution

$$\text{Cat}(y \mid \mathbf{C} = \begin{bmatrix} c_0 = 0.1 \\ c_1 = 0.7 \\ c_2 = 0.2 \end{bmatrix})$$

Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

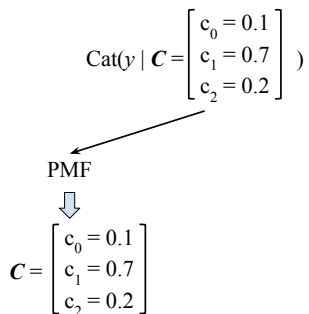


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

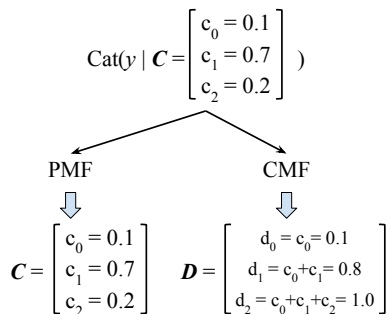


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

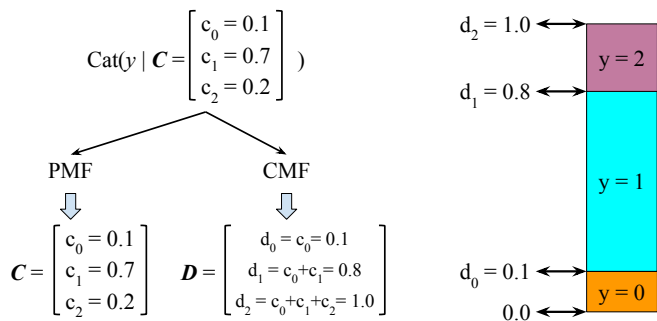


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

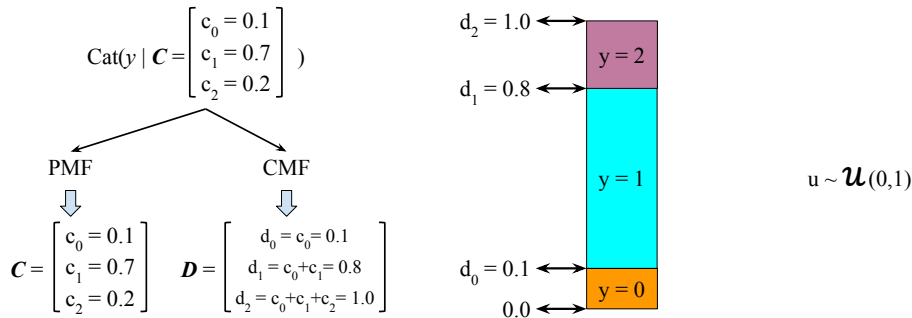


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

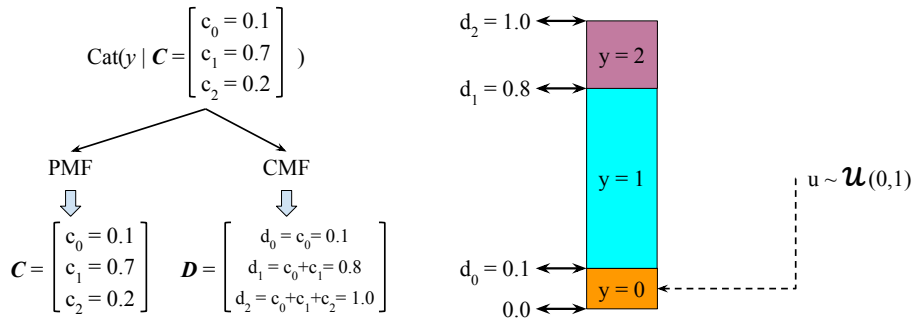


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

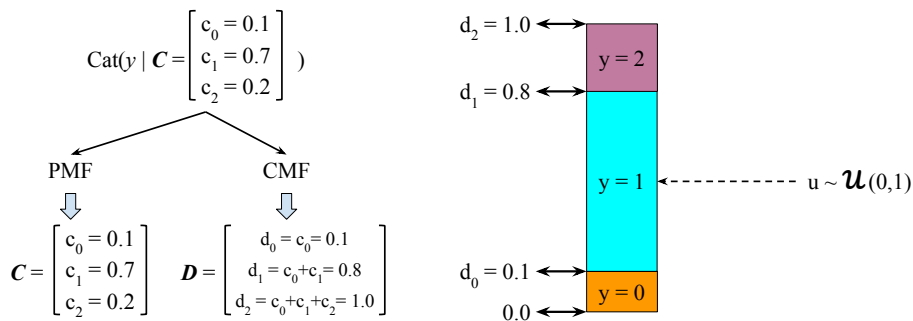


Figure: Sampling a categorical distribution using a Uniform sampler

Sampling a Categorical Distribution

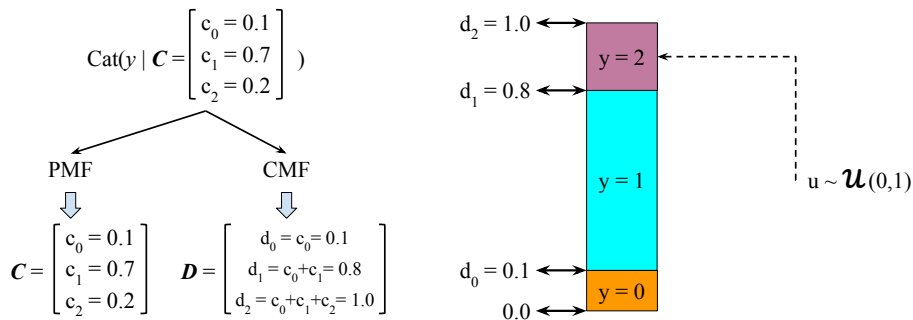


Figure: Sampling a categorical distribution using a Uniform sampler

Subsection 2

Deep Autoregressive Models

Model Specification

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 $\mathbf{x} \in \mathbb{R}^{784}$.

Model Specification

Assume we just have MINST 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.

Model Specification

Assume we just have MINST 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.

Model Specification

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 $x_i \in \{0, 1, 2, \dots, 255\}$

Model Specification

Assume we just have MINST 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:

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

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

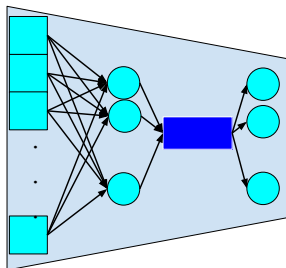
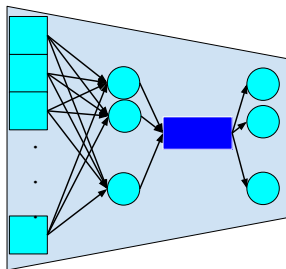


Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

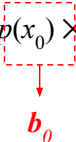


$$\mathbf{W}_d, \mathbf{b}_d$$

Figure: Using logistic regression for generative modeling

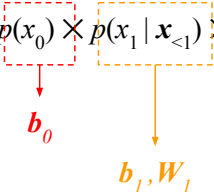
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$


The diagram shows the term $p(x_0)$ in the joint probability distribution formula. A red dashed box is drawn around $p(x_0)$, and a red arrow points from the box down to the parameter b_0 .

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$


The diagram illustrates the decomposition of a joint probability distribution $p(\mathbf{x})$ into a product of conditional distributions. The first term, $p(x_0)$, is highlighted with a red dashed box and a red arrow pointing to the parameter vector \mathbf{b}_0 . The second term, $p(x_1 | \mathbf{x}_{<1})$, is highlighted with an orange dashed box and an orange arrow pointing to the parameters \mathbf{b}_1 and \mathbf{W}_1 .

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

\mathbf{b}_0

$\mathbf{b}_1, \mathbf{W}_1$

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

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

\mathbf{b}_0

$\mathbf{b}_1, \mathbf{W}_1$

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

$\mathbf{W}_{D-1}, \mathbf{b}_{D-1}$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

\mathbf{b}_0
 $\mathbf{b}_1, \mathbf{W}_1$
 $\mathbf{b}_d, \mathbf{W}_d$
 $\mathbf{W}_{D-1}, \mathbf{b}_{D-1}$

$$x_d \in \{0, 1, \dots, 255\} \Rightarrow \begin{cases} \mathbf{b}_d \in \mathbb{R}^{256} \\ \mathbf{W}_d \in \mathbb{R}^{256 \times d} \end{cases} \quad \forall \quad 0 \leq d \leq D-1$$

Figure: Using logistic regression for generative modeling

$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times p(x_d | \mathbf{x}_{<d}) \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$

\mathbf{b}_0
 $\mathbf{b}_1, \mathbf{W}_1$
 $\mathbf{b}_d, \mathbf{W}_d$
 $\mathbf{W}_{D-1}, \mathbf{b}_{D-1}$

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

Figure: Using logistic regression for generative modeling

Distance Metric

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

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

Distance Metric

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

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

Distance Metric

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

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

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_{data} and p_{θ} , thus we can use KL divergence as:

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

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 θ , 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:

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 $p_{\text{data}}(\mathbb{X})$

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_{data}

Monte Carlo Estimation

Consider the following expectation:

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

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 $\mathbf{x}_1, \dots, \mathbf{x}_N$.

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:

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

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_{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_{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:

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

Parametric Density Calculation

$$\theta = \{ \textcolor{red}{b}_0, \textcolor{brown}{b}_1, \textcolor{brown}{W}_1, \dots, \textcolor{blue}{b}_d, \textcolor{blue}{W}_d, \dots, \textcolor{violet}{b}_{D-1}, \textcolor{violet}{W}_{D-1} \}$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \textcolor{red}{b}_0, \textcolor{orange}{b}_1, \textcolor{brown}{W}_1, \dots, \textcolor{blue}{b}_d, \textcolor{blue}{W}_d, \dots, \textcolor{violet}{b}_{D-1}, \textcolor{violet}{W}_{D-1} \}$$

x

$x_0=255$

$x_1=126$

\vdots

$x_{d-1}=65$

$x_d=23$

\vdots

$x_{D-1}=0$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \textcolor{red}{b}_0, \textcolor{brown}{b}_1, \textcolor{brown}{W}_1, \dots, \textcolor{blue}{b}_d, \textcolor{blue}{W}_d, \dots, \textcolor{violet}{b}_{D-1}, \textcolor{violet}{W}_{D-1} \}$$

\mathbf{x}

$x_0=255$

$x_1=126$

\vdots

$x_{d-1}=65$

$x_d=23$

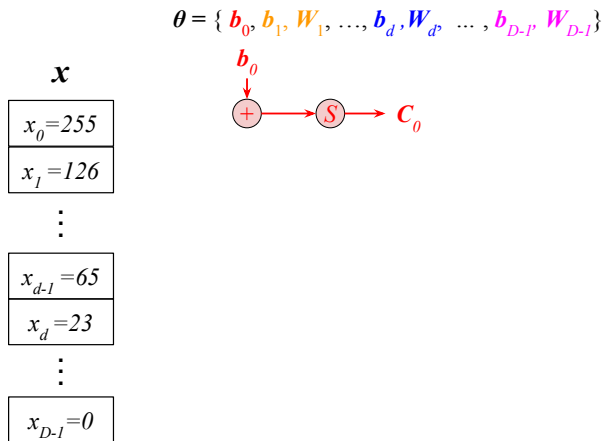
\vdots

$x_{D-1}=0$

$$p(\mathbf{x}) = p(x_0)p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

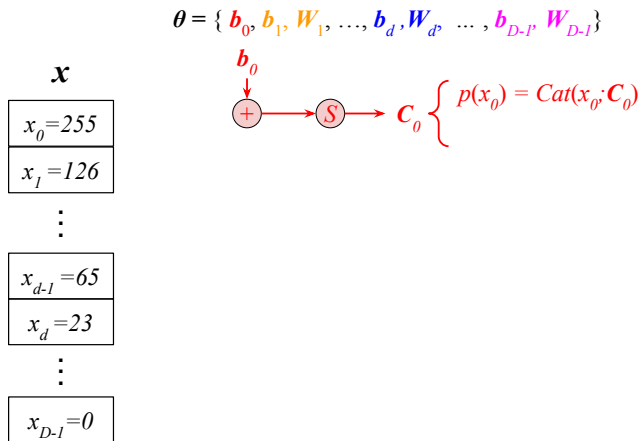
Parametric Density Calculation



$$p(\mathbf{x}) = p(x_0)p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

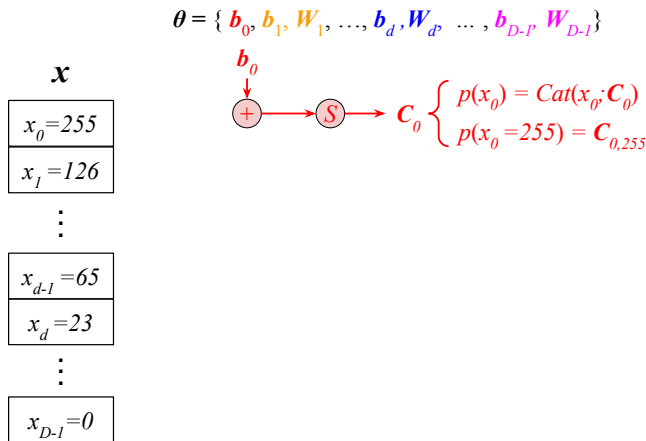
Parametric Density Calculation



$$p(\mathbf{x}) = p(x_0)p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

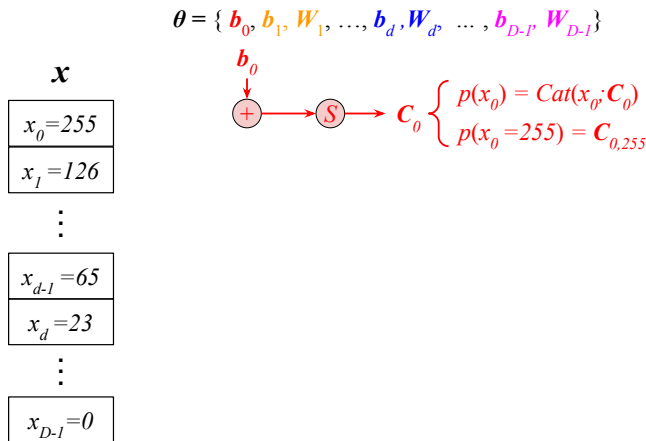
Parametric Density Calculation



$$p(\mathbf{x}) = p(x_0)p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

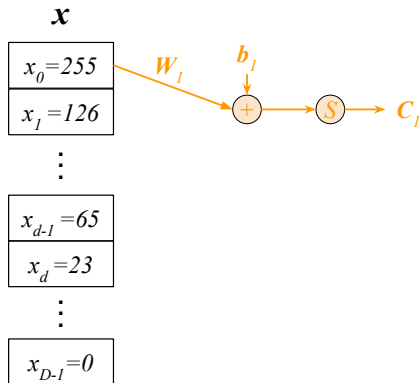


$$p(\mathbf{x}) = C_{0,255} p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

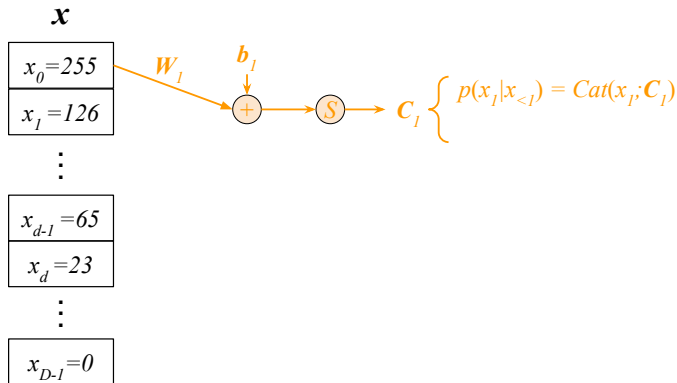


$$p(\mathbf{x}) = \mathbf{C}_{0,255} p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

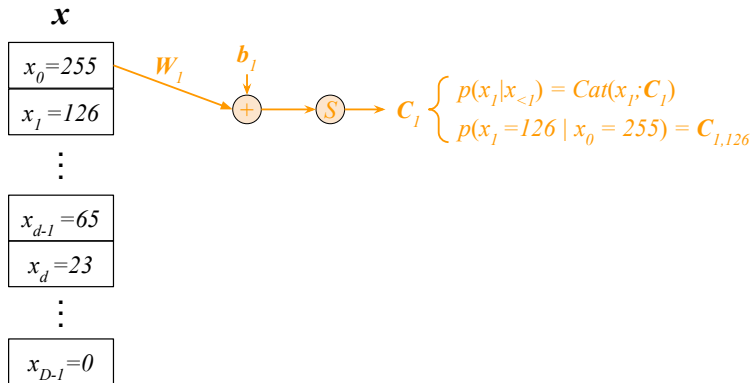


$$p(\mathbf{x}) = \mathbf{C}_{0,255} p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$



$$p(\mathbf{x}) = \mathbf{C}_{0,255} p(x_1 | \mathbf{x}_{<1}) \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

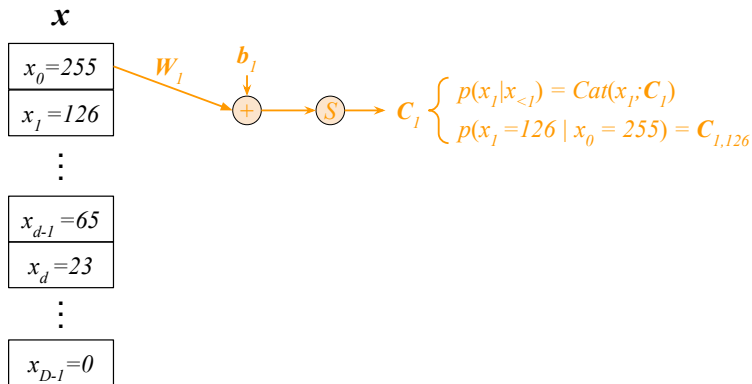
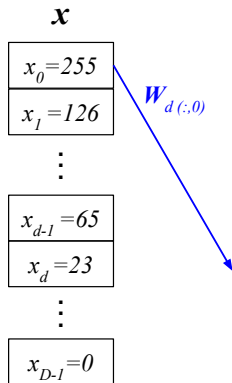


Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \textcolor{red}{b}_0, \textcolor{orange}{b}_1, \textcolor{brown}{W}_1, \dots, \textcolor{blue}{b}_d, \textcolor{blue}{W}_d, \dots, \textcolor{violet}{b}_{D-1}, \textcolor{violet}{W}_{D-1} \}$$

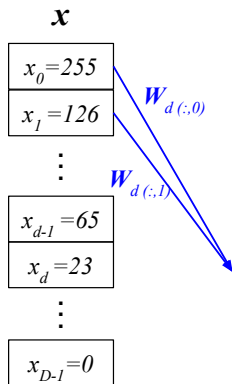


$$p(\mathbf{x}) = \textcolor{red}{C}_{0,255} \textcolor{orange}{C}_{1,126} \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

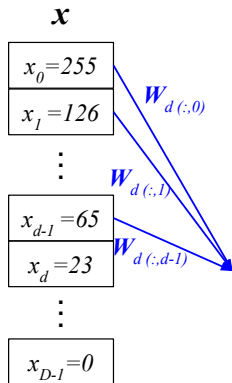


$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

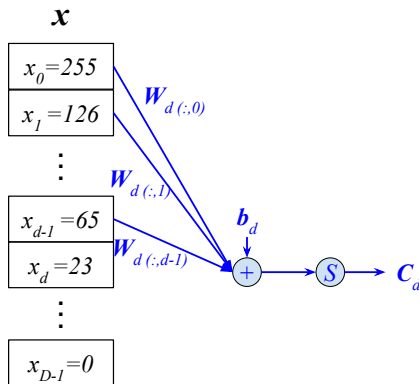


$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

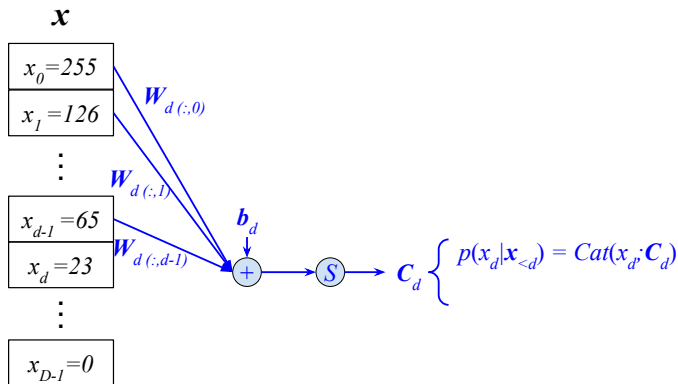


$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$



$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots p(x_d | \mathbf{x}_{<d}) \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

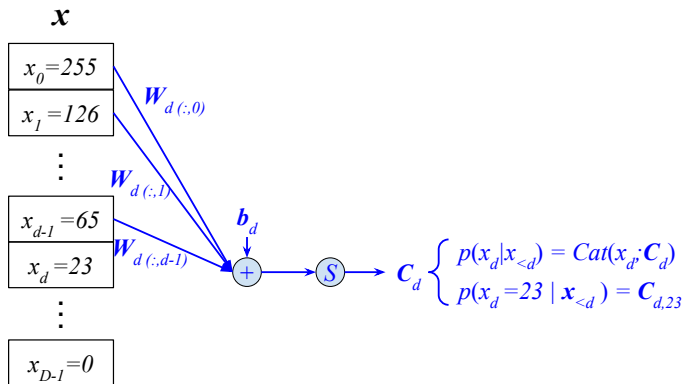
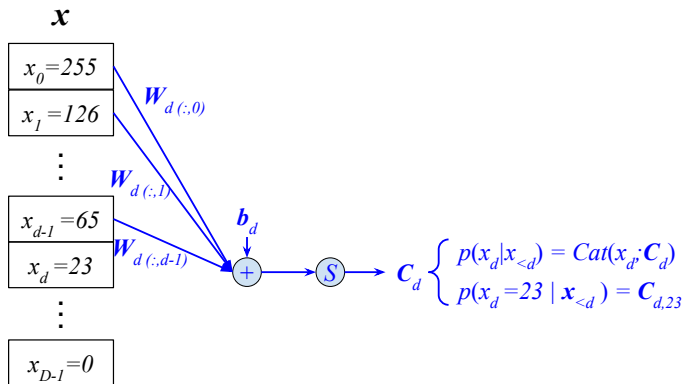


Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

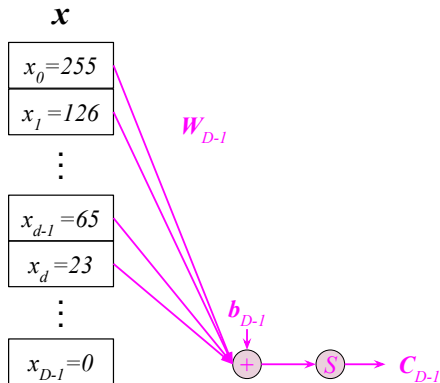


$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots \mathbf{C}_{d,23} \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

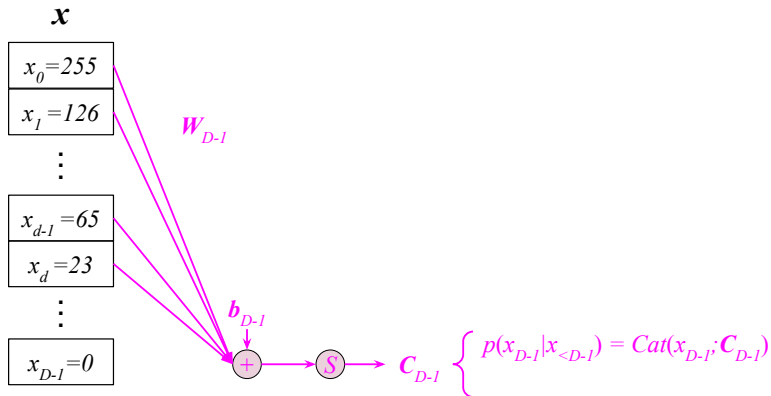


$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots \mathbf{C}_{d,23} \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

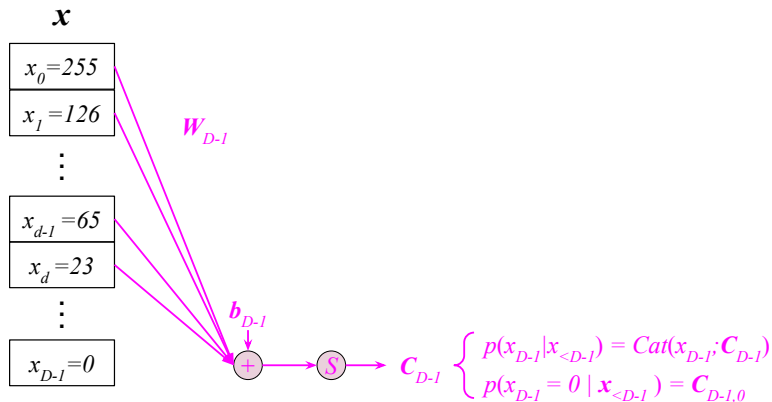


$$p(\mathbf{x}) = C_{0,255} C_{1,126} \dots C_{d,23} \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$



$$p(\mathbf{x}) = \mathbf{C}_{0,255} \mathbf{C}_{1,126} \dots \mathbf{C}_{d,23} \dots p(x_{D-1} | \mathbf{x}_{<D-1})$$

Figure: Calculating the likelihood as a function of model parameters

Parametric Density Calculation

$$\theta = \{ \mathbf{b}_0, \mathbf{b}_1, \mathbf{W}_1, \dots, \mathbf{b}_d, \mathbf{W}_d, \dots, \mathbf{b}_{D-1}, \mathbf{W}_{D-1} \}$$

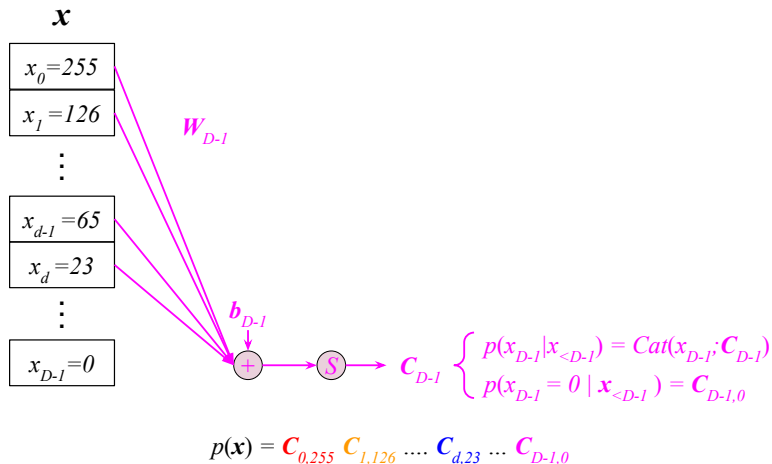


Figure: Calculating the likelihood as a function of model parameters

Sampling from a Generative Model

$$\theta^{\star} = \{ \textcolor{red}{b}_0^{\star}, \textcolor{orange}{b}_1^{\star}, \textcolor{orange}{w}_1^{\star}, \dots, \textcolor{blue}{b}_d^{\star}, \textcolor{blue}{w}_{d'}^{\star} \dots, \textcolor{violet}{b}_{D-P}^{\star}, \textcolor{violet}{w}_{D-I}^{\star} \}$$

Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^{\star} = \{ \textcolor{red}{b}_0^{\star}, \textcolor{orange}{b}_1^{\star}, \textcolor{brown}{w}_1^{\star}, \dots, \textcolor{blue}{b}_d^{\star}, \textcolor{blue}{w}_{d'}^{\star} \dots, \textcolor{violet}{b}_{D-I}^{\star}, \textcolor{violet}{w}_{D-I}^{\star} \}$$

x

$x_0 = ?$

$x_1 = ?$

\vdots

$x_{d-1} = ?$

$x_d = ?$

\vdots

$x_{D-1} = ?$

Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

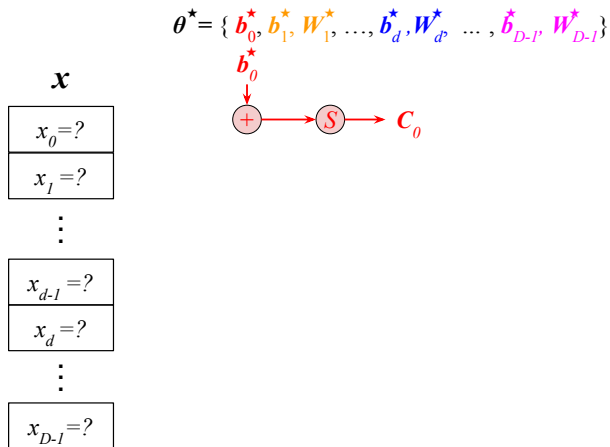


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

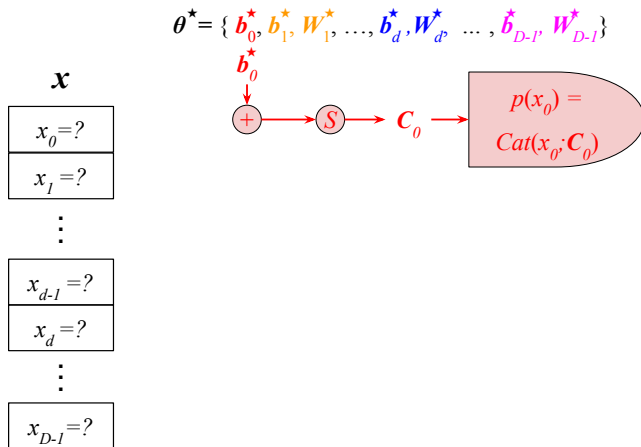


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

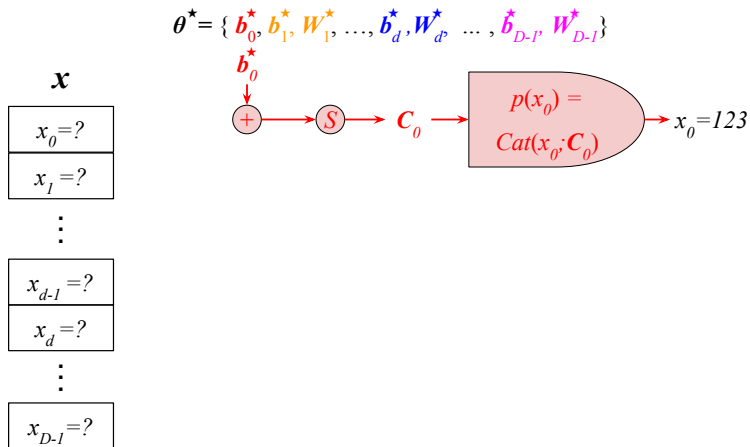


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

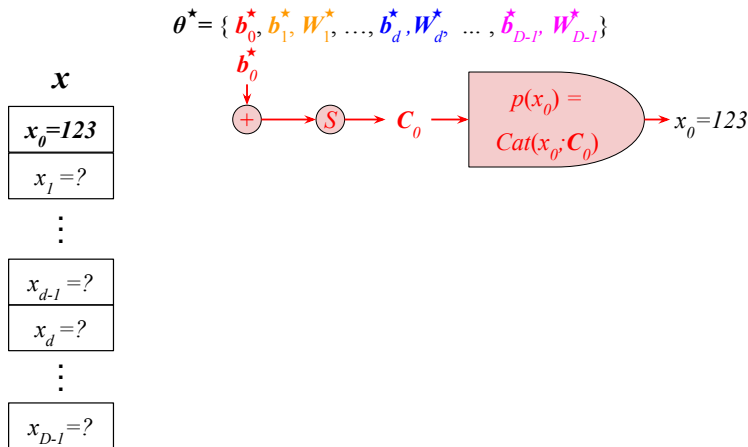


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^{\star} = \{ \textcolor{red}{b}_0^{\star}, \textcolor{brown}{b}_1^{\star}, \textcolor{brown}{w}_1^{\star}, \dots, \textcolor{blue}{b}_d^{\star}, \textcolor{blue}{w}_{d'}^{\star} \dots, \textcolor{violet}{b}_{D-I}^{\star}, \textcolor{violet}{w}_{D-I}^{\star} \}$$

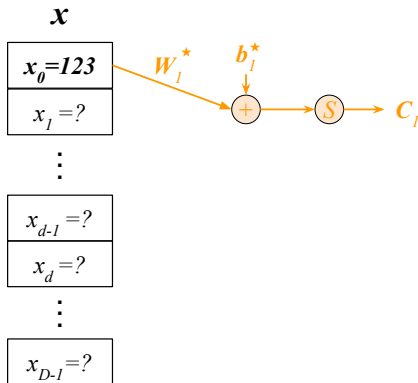


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

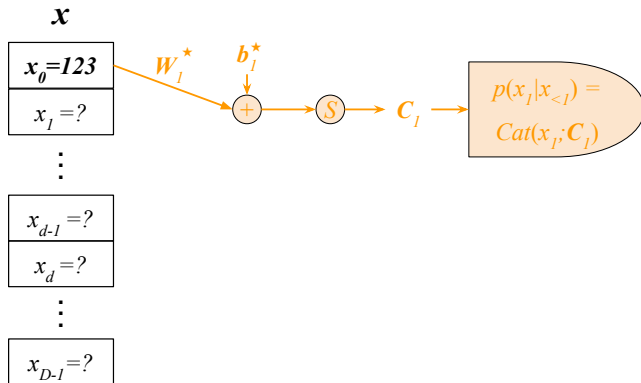


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

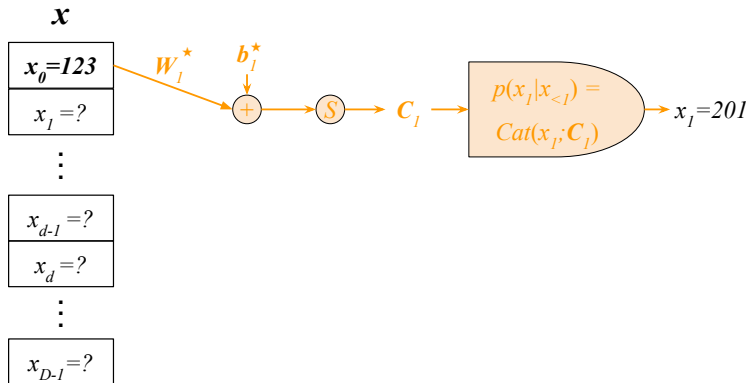


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

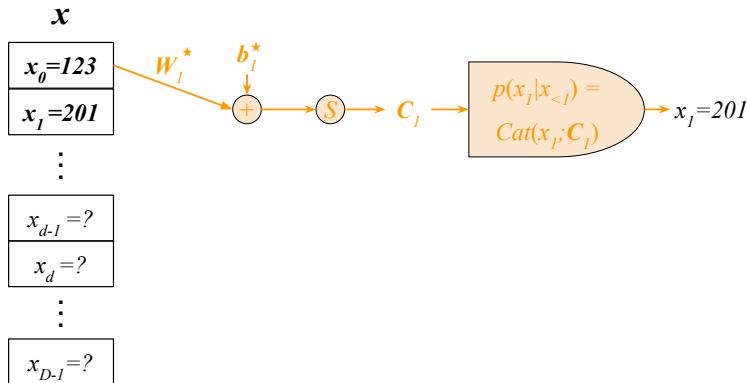


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^{\star} = \{ \mathbf{b}_0^{\star}, \mathbf{b}_1^{\star}, \mathbf{W}_1^{\star}, \dots, \mathbf{b}_d^{\star}, \mathbf{W}_d^{\star}, \dots, \mathbf{b}_{D-1}^{\star}, \mathbf{W}_{D-1}^{\star} \}$$

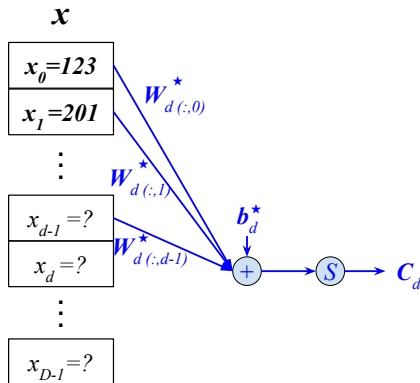


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_{d'}^* \dots, \mathbf{b}_{D-I}^*, \mathbf{W}_{D-I}^* \}$$

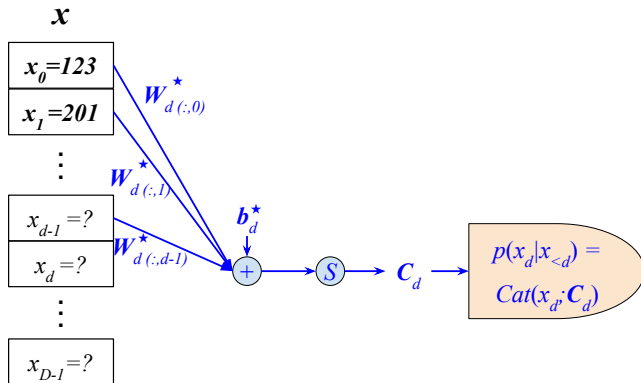


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^*, \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

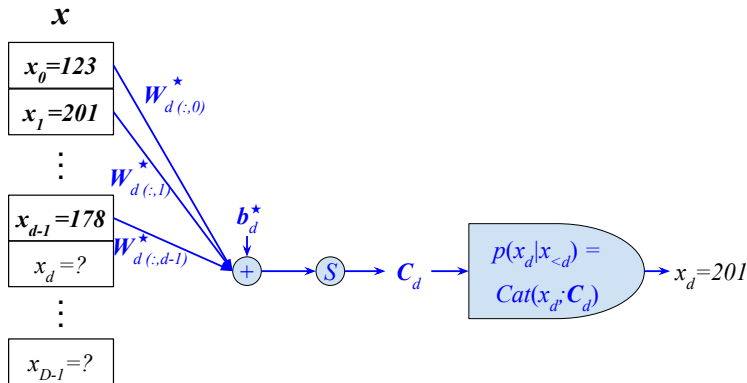


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^*, \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

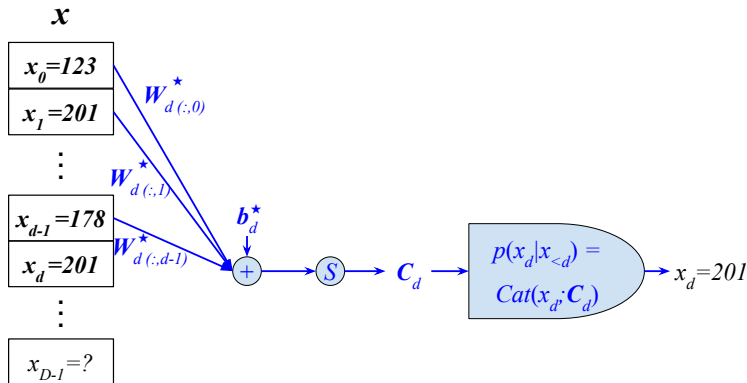


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^{\star} = \{ \textcolor{red}{b}_0^{\star}, \textcolor{orange}{b}_1^{\star}, \textcolor{orange}{w}_1^{\star}, \dots, \textcolor{blue}{b}_d^{\star}, \textcolor{blue}{w}_d^{\star} \dots, \textcolor{violet}{b}_{D-1}^{\star}, \textcolor{violet}{w}_{D-1}^{\star} \}$$

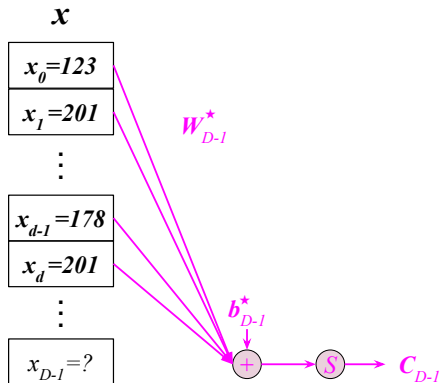


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

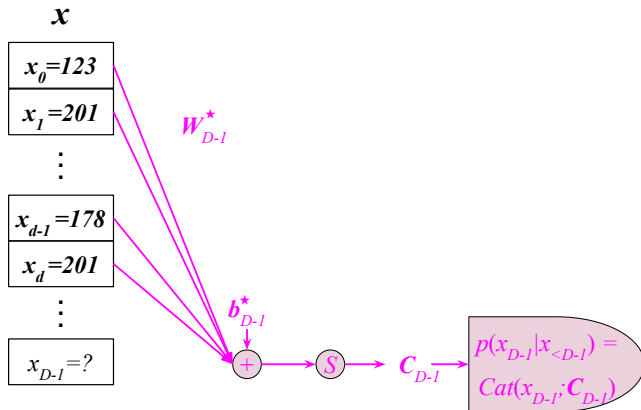


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

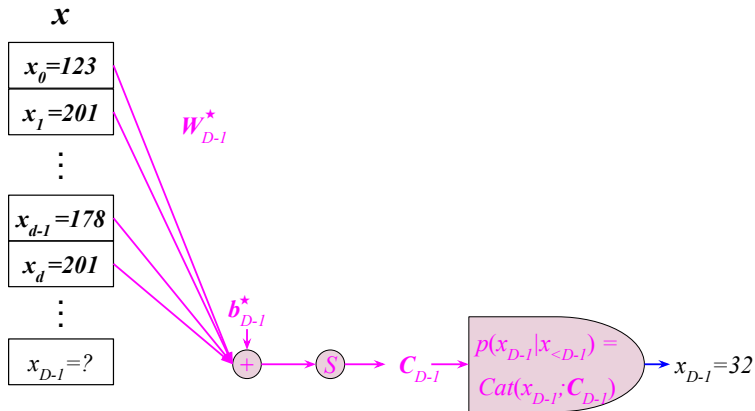


Figure: Sampling a trained Autoregressive Model

Sampling from a Generative Model

$$\theta^* = \{ \mathbf{b}_0^*, \mathbf{b}_1^*, \mathbf{W}_1^*, \dots, \mathbf{b}_d^*, \mathbf{W}_d^* \dots, \mathbf{b}_{D-1}^*, \mathbf{W}_{D-1}^* \}$$

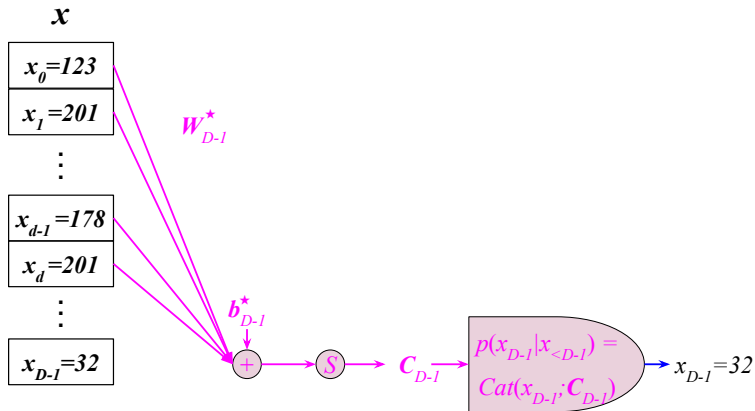


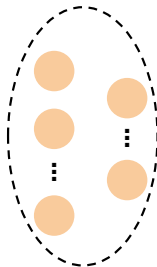
Figure: Sampling a trained Autoregressive Model

Section 7

Extensions

Some of Autoregressive Modeling Extensions

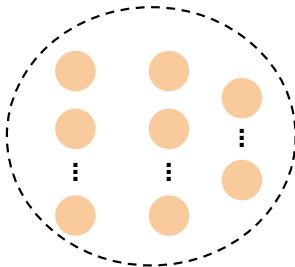
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$



Fully Visible Sigmoid Belief Networks
(FVSN)

Some of Autoregressive Modeling Extensions

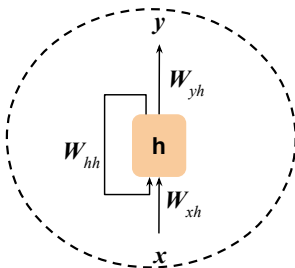
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$



Neural Autoregressive Density Estimation
(NADE)

Some of Autoregressive Modeling Extensions

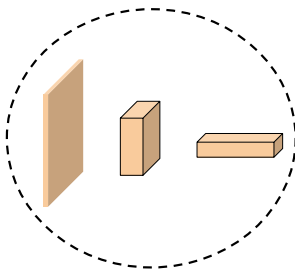
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$



Pixel Recurrent Neural Networks
(PixelRNN)

Some of Autoregressive Modeling Extensions

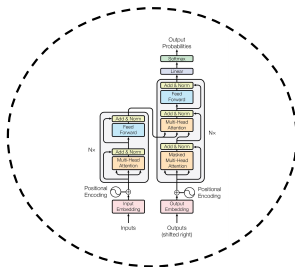
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$



Pixel Convolutional Neural Networks
(PixelCNN)

Some of Autoregressive Modeling Extensions

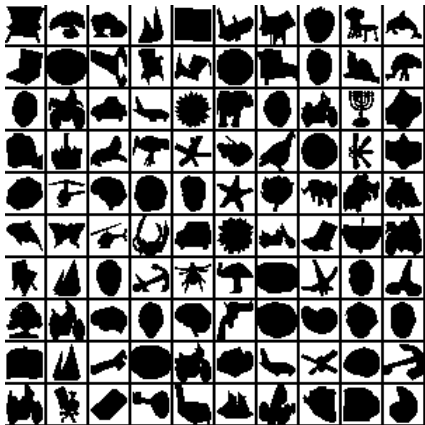
$$p(\mathbf{x}) = p(x_0) \times p(x_1 | \mathbf{x}_{<1}) \times \dots \times \boxed{p(x_d | \mathbf{x}_{<d})} \times \dots \times p(x_{D-1} | \mathbf{x}_{<D-1})$$



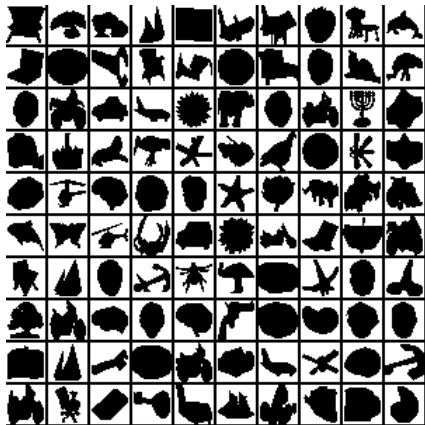
Transformer
(ChatGPT)

Section 8

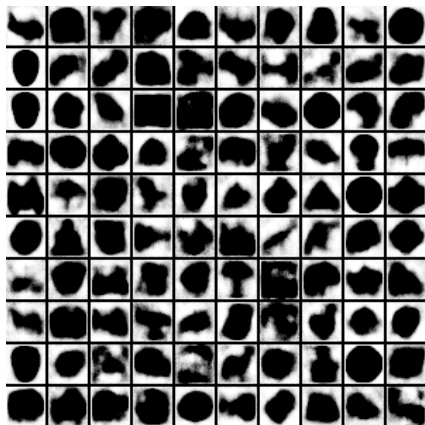
Results



(a) Dataset samples



(a) Dataset samples



(b) Generated samples

Figure: FVSBN performance over Caltech 101 dataset (source: [5])



Figure: NADE performance over BMNIST dataset (source: [6])

occluded

completions

original

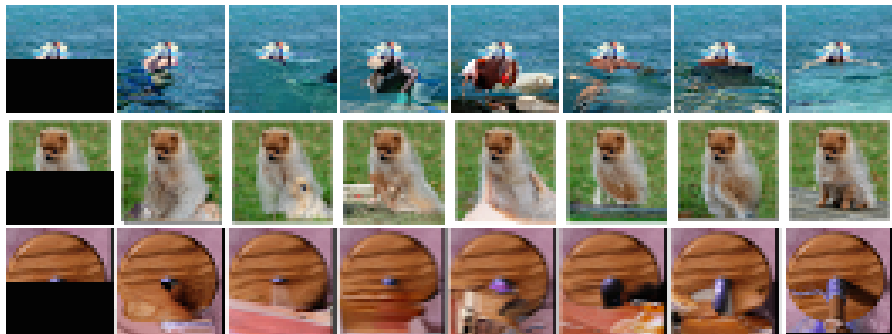


Figure: Pixel RNN results in image completion (source: [7])

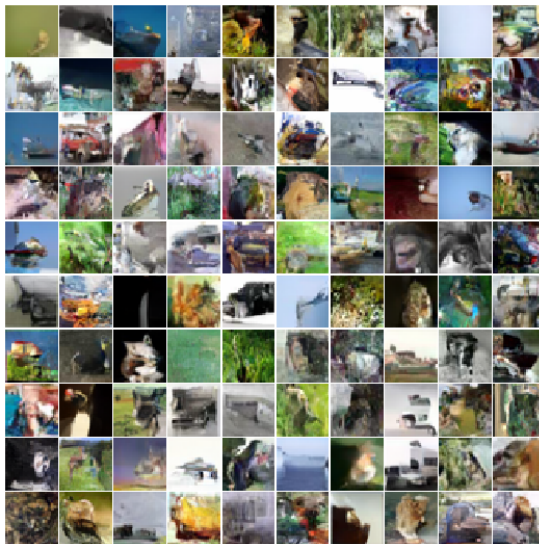


Figure: Samples from our PixelCNN model trained on CIFAR-10 (source: [8])

Section 9

Applications

Adversarial Robustness

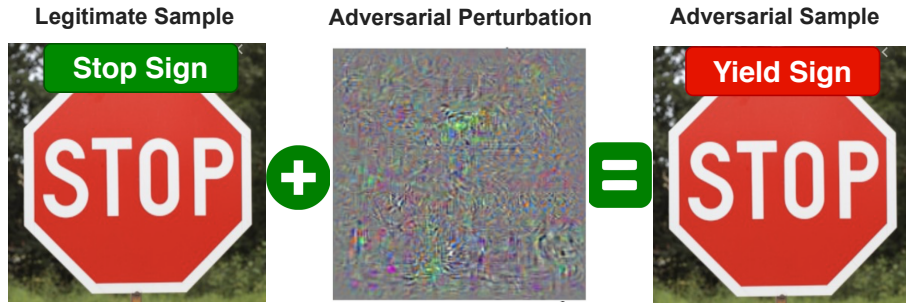
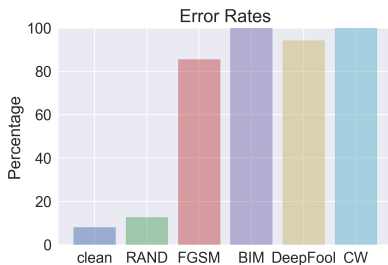


Figure: Different adversarial attacks to Frog image from Cifar10 dataset (source: [9])

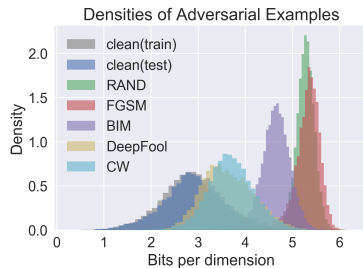
Adversarial Robustness



Figure: Sample adversarial attack to deep learning architectures (source: [9])



(a) Error rate



(b) Density change

Figure: Using autoregressive models to detect adversarial samples (source: [9])

Thank You!

Thank you for your attention!

Do you have any questions or comments?

Contact Information

Sajjad Amini
Email: 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, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al.,
“Photorealistic text-to-image diffusion models with deep language understanding,”

Advances in Neural Information Processing Systems, vol. 35, pp. 36479–36494, 2022.



Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi,

“Palette: Image-to-image diffusion models,”
in *ACM SIGGRAPH 2022 Conference Proceedings*, 2022, pp. 1–10.



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.

References II



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.



Aäron Van Den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu,
“Pixel recurrent neural networks,”
in International conference on machine learning. PMLR, 2016, pp. 1747–1756.



Tim Salimans, Andrej Karpathy, Xi Chen, and Diederik P Kingma,
“Pixelcnn++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications,”
arXiv preprint arXiv:1701.05517, 2017.



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.