

Deep Generative Models

3. Autoregressive Models



- 국가수리과학연구소 산업수학혁신센터 김민중

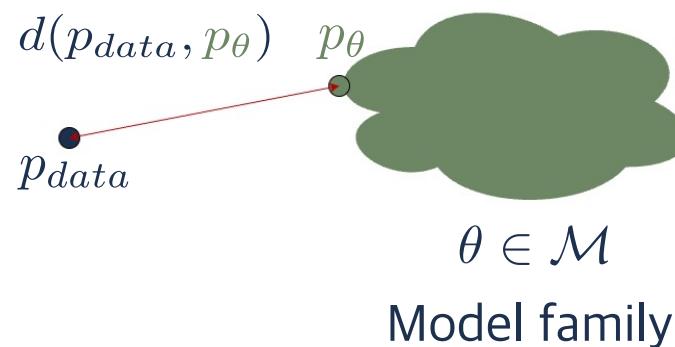
Learning a generative model

- We are given a training dataset of examples.



$$\begin{aligned} \mathbf{x}_i &\sim p_{data} \\ i &= 1, 2, \dots, N \end{aligned}$$

- Generation: sample \mathbf{x}_{new} should look like training set(sampling)
- Density estimation
- Unsupervised representation learning: learn what these images have in common features
- 1st question: How to represent probability distribution?
- 2nd question: How to learn it



Recap: Bayesian networks vs neural models

- Using Chain rule, full general form is

$$p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2|x_1)p(x_3|x_1, x_2)p(x_4|x_1, x_2, x_3)$$

- Bayes Net under conditional independencies

$$p(x_1, x_2, x_3, x_4) \approx p(x_1)p(x_2|x_1)p(x_3|x_1, x_2)p(x_4|x_1, x_2, x_3)$$

- Neural Models

$$p(x_1, x_2, x_3, x_4)$$

$$\approx p(x_1)p_{\text{neural}}(x_2|x_1)p_{\text{neural}}(x_3|x_1, x_2)p_{\text{neural}}(x_4|x_1, x_2, x_3)$$

- Assume specific functional form for the conditionals
- A sufficiently deep neural net can approximate any function

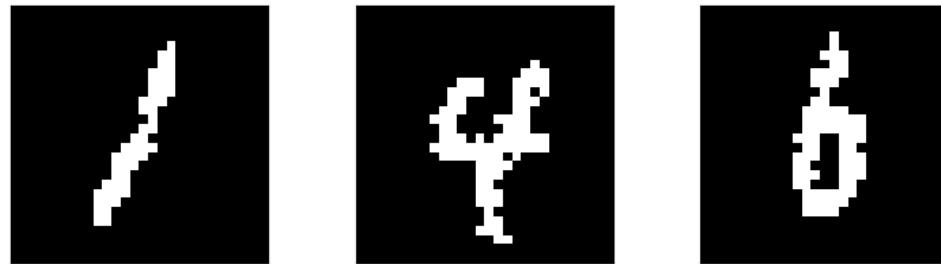
Recap: Neural Models for classification

- For classification, we care about $p(Y = 1|x)$ and assume that
$$p(Y = 1|x; \alpha) = f(x, \alpha)$$
- Linear dependence:
 - Let $z(\alpha, x) := \alpha_0 + \sum_{i=1}^d \alpha_i x_i$
 - $p(Y = 1|x; \alpha) = \sigma(z(\alpha, x))$ where $\sigma(z) = 1/(1 + e^{-z})$
- Non-linear dependence: let $h(W, b, x) = f(Wx + b)$ be a non-linear transformation of the inputs (features)

$$\begin{aligned} p_{neural}(Y = 1|x; \alpha, W, b) &= \sigma\left(z(\alpha, h(W, b, x))\right) \\ &= \sigma\left(\alpha_0 + \sum_{i=1}^h \alpha_i h_i\right) \end{aligned}$$

- Parameters: α, W, b
- Can repeat multiple times to get a neural network

Motivating Example: MNIST



- Given: a dataset D of handwritten digits (binarized MNIST)
- Each image has $d = 28 \times 28 = 784$ pixels. Each pixel can either be black (0) or white (1)
- **Goal:** Learn a probability distribution $p(\mathbf{x}) = p(x_1, \dots, x_{784})$ over $x_i \in \{0,1\}^{784}$ such that when $\mathbf{x} \sim p_\theta(\mathbf{x})$, \mathbf{x} looks like a digit
- Two step process:
 - Parameterize a model family $\{p_\theta(\mathbf{x}), \theta \in \mathcal{M}\}$
 - Search for model parameters θ based on training data D

Autoregressive model

- Given a dataset D of d -dimensional datapoint \mathbf{x}
- For simplicity, let $\mathbf{x} \in \{0,1\}^d$
- By the chain rule, (no conditional independence)

$$p(\mathbf{x}) = p(x_1) \prod_{i=2}^d p(x_i | \mathbf{x}_{<i})$$

- where $\mathbf{x}_{<i} = (x_1, x_2, \dots, x_{i-1})^T$ denotes the vector of random variables with index less than i
- Such a Bayesian network is said to obey the autoregressive property

Autoregressive model

- Fix an ordering of variables
- x_i random variable depends on x_1, x_2, \dots, x_{i-1}
- To fully specify $p(x_i|x_{<i})$, we need to specify a probability for 2^{i-1} configurations of the variable x_1, x_2, \dots, x_{i-1}

Autoregressive model

- The conditionals are specified as parameter functions with a fixed number of parameters
- Assume $p(x_i | \mathbf{x}_{<i}) \sim \text{Bernoulli}$ random variable and learn a function mapping $\mathbf{x}_{<i}$ to the mean of this Bernoulli distribution
- Hence, for $i > 1$,

$$p_{\theta_i}(x_i | \mathbf{x}_{<i}) = \text{Bern}(x_i | f_i(\mathbf{x}_{<i}))$$

- where θ_i denotes the set of parameters used to specify the mean function $f_i: \{0,1\}^{i-1} \rightarrow (0,1)$
- # of parameters of an autoregressive generative model is $\sum_{i=1}^d |\theta_i|$

Fully-visible sigmoid belief network(FVSBN)

- Specify the function as a linear combination of the input elements followed by a sigmoid non-linearity

$$p_{\theta_i}(x_i | \mathbf{x}_{<i}) = \text{Bern}(x_i | f_i(\mathbf{x}_{<i}))$$

$$f_i(\mathbf{x}_{<i}) = \sigma(\mathbf{x}_{<i}^T \mathbf{w}^{(i)} + b_i)$$

- where σ is a sigmoid, $\mathbf{w}^{(i)} \in \mathbb{R}^{i-1}$ and $b_i \in \mathbb{R}$. i.e., $\theta_i = \{\mathbf{w}^{(i)}, b_i\}$
- Since f_i requires i parameters, the total number of parameters is $\sum_{i=1}^d i = O(d^2)$

Fully-visible sigmoid belief network(FVSBN)

- E.g., $\mathbf{x} = (x_1, x_2, x_3)^T$ with $x_i \in \{0,1\}$, then

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

- FVSBN model is

$$p_{\theta}(\mathbf{x}) = p_{\theta_1}(x_1)p_{\theta_2}(x_2|x_1)p_{\theta_3}(x_3|x_1, x_2)$$

- where

$$p_{\theta_1}(x_1) = \text{Bern}(x_1|\theta_1), 0 < \theta_1 < 1,$$

$$p_{\theta_2}(x_2|x_1) = \text{Bern}(x_2|f_2(x_1)), f_2(x_1) = \sigma(w_1^{(2)}x_1 + b_2)$$

$$p_{\theta_3}(x_3|x_1, x_2) = \text{Bern}(x_3|f_3(x_1, x_2)),$$

$$f_3(x_1, x_2) = \sigma(w_1^{(3)}x_1 + w_2^{(3)}x_2 + b_3)$$

- # of trainable parameters is $\sum_{i=1}^3 |\theta_i| = 1 + 2 + 3$

Fully-visible sigmoid belief network(FVSBN)

- FVSBN model is

$$p_{\theta}(\mathbf{x}) = p_{\theta_1}(x_1)p_{\theta_2}(x_2|x_1)p_{\theta_3}(x_3|x_1, x_2)$$

- where

$$p_{\theta_1}(x_1) = \text{Bern}(x_1|\theta_1), 0 < \theta_1 < 1,$$

$$p_{\theta_2}(x_2|x_1) = \text{Bern}(x_2|f_2(x_1)), f_2(x_1) = \sigma(w_1^{(2)}x_1 + b_2)$$

$$p_{\theta_3}(x_3|x_1, x_2) = \text{Bern}(x_3|f_3(x_1, x_2)),$$

$$f_3(x_1, x_2) = \sigma(w_1^{(3)}x_1 + w_2^{(3)}x_2 + b_3)$$

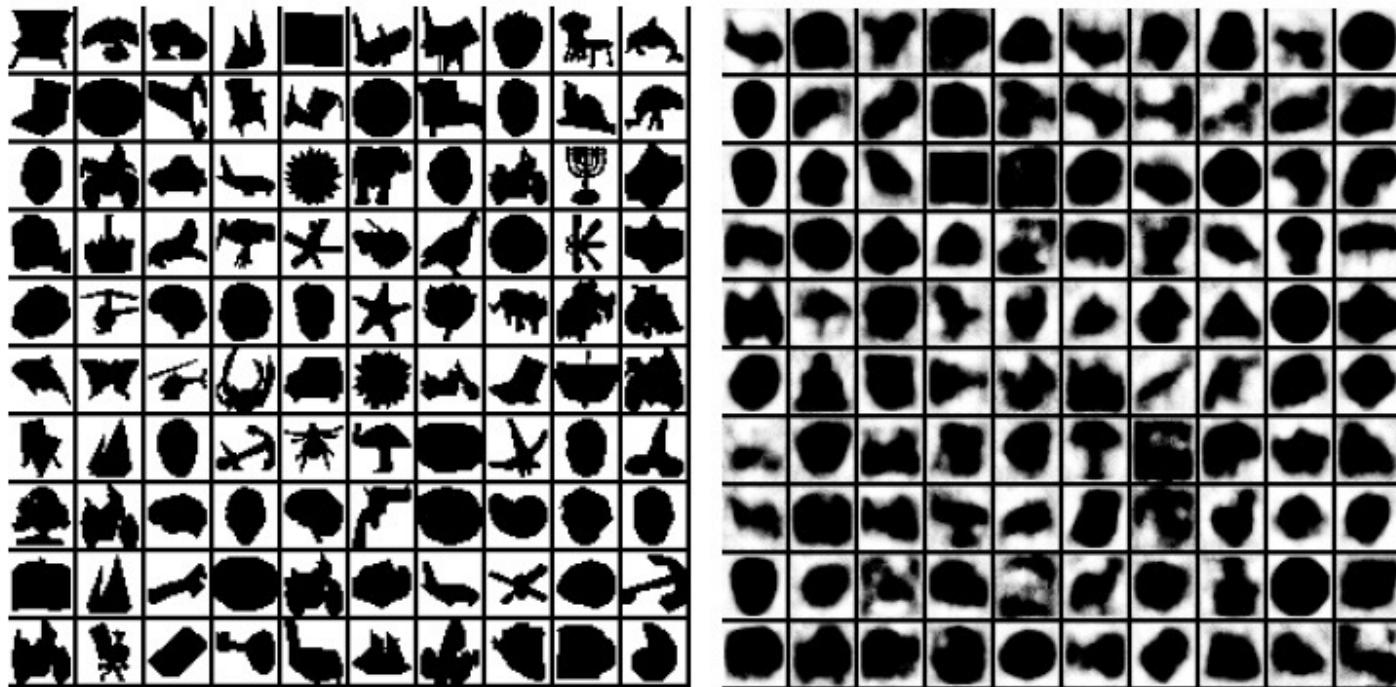
- # of trainable parameters is $\sum_{i=1}^3 |\theta_i| = 1 + 2 + 3$

- Since $\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{(1-x)}$,

$$\begin{aligned} p_{\theta}(\mathbf{x}) \\ = \theta_1^{x_1}(1 - \theta_1)^{(1-x_1)}f_2(x_1)^{x_2}(1 - f_2(x_1))^{(1-x_2)}f_3(x_1, x_2)^{x_3}(1 \\ - f_3(x_1, x_2))^{(1-x_3)} \end{aligned}$$

Fully-visible sigmoid belief network(FVSBN)

- Training data on the left (Caltech 101 Silhouettes)
- Samples from the model on the right.
- Figure from Learning Deep Sigmoid Belief Networks with Data Augmentation, 2015.



Fully-visible sigmoid belief network(FVSBN)

- More flexible parametrizations for the mean function e.g., MLP
- For example, consider 1 hidden layer neural network. i.e.,

$$\begin{aligned}\mathbf{h}_i(\mathbf{x}_{*)} &= \sigma(A_i \mathbf{x}_{*} + \mathbf{c}_i) \\ p_{\theta_i}(x_i | \mathbf{x}_{})) \\ f_i(\mathbf{x}_{*}) &= \sigma(\mathbf{h}_i^T \mathbf{w}^{(i)} + b_i)\end{aligned}***$$

- where $\mathbf{h}_i \in \mathbb{R}^h$ denotes the hidden layer activations for MLP and $\theta_i \in \{A_i \in \mathbb{R}^{h \times (i-1)}, \mathbf{c}_i \in \mathbb{R}^h, \mathbf{w}^{(i)} \in \mathbb{R}^h, b_i \in \mathbb{R}\}$ are the set of parameters for the mean function
- # of parameters is dominated by $|A_i|$ so given by $O(d^2 h)$

$$\mathbf{h}_2 = \sigma \left(\underbrace{\begin{pmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{pmatrix}}_{A_2} x_1 + \begin{pmatrix} \cdot \\ \cdot \\ \cdot \\ \cdot \end{pmatrix} \right), \mathbf{h}_3 = \sigma \left(\underbrace{\begin{pmatrix} \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{pmatrix}}_{A_3} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} + \begin{pmatrix} \cdot \\ \cdot \\ \cdot \end{pmatrix} \right)$$

Neural Autoregressive Density estimator(NADE)

- An alternation of MLP-based parametrization
- In NADE, parameters are shared across the functions used for evaluating the conditionals
- In particular,

$$\begin{aligned} \mathbf{h}_i(\mathbf{x}_{<i}) &= \sigma(A_{.,<i}\mathbf{x}_{<i} + \mathbf{c}) \\ p_{\theta_i}(x_i|\mathbf{x}_{<i}) &= \text{Bern}(x_i|f_i(\mathbf{x}_{<i})) \\ f_i(\mathbf{x}_{<i}) &= \sigma(\mathbf{h}_i^T \mathbf{w}^{(i)} + b_i) \end{aligned}$$

- where $\theta \in \{A \in \mathbb{R}^{h \times d}, \mathbf{c} \in \mathbb{R}^h, \{\mathbf{w}^{(i)} \in \mathbb{R}^h\}_{i=1}^d, \{b_i \in \mathbb{R}\}_{i=1}^d\}$ is the full parameters for the mean function
- The weight matrix A and the bias vector \mathbf{c} are shared across the conditionals
- # of parameters gets reduced from $O(d^2h)$ to $O(dh)$

Neural Autoregressive Density estimator(NADE)

- Samples from a model trained on MNIST on the left
- Figure from The Neural Autoregressive Distribution Estimator, 2011



Thanks
