

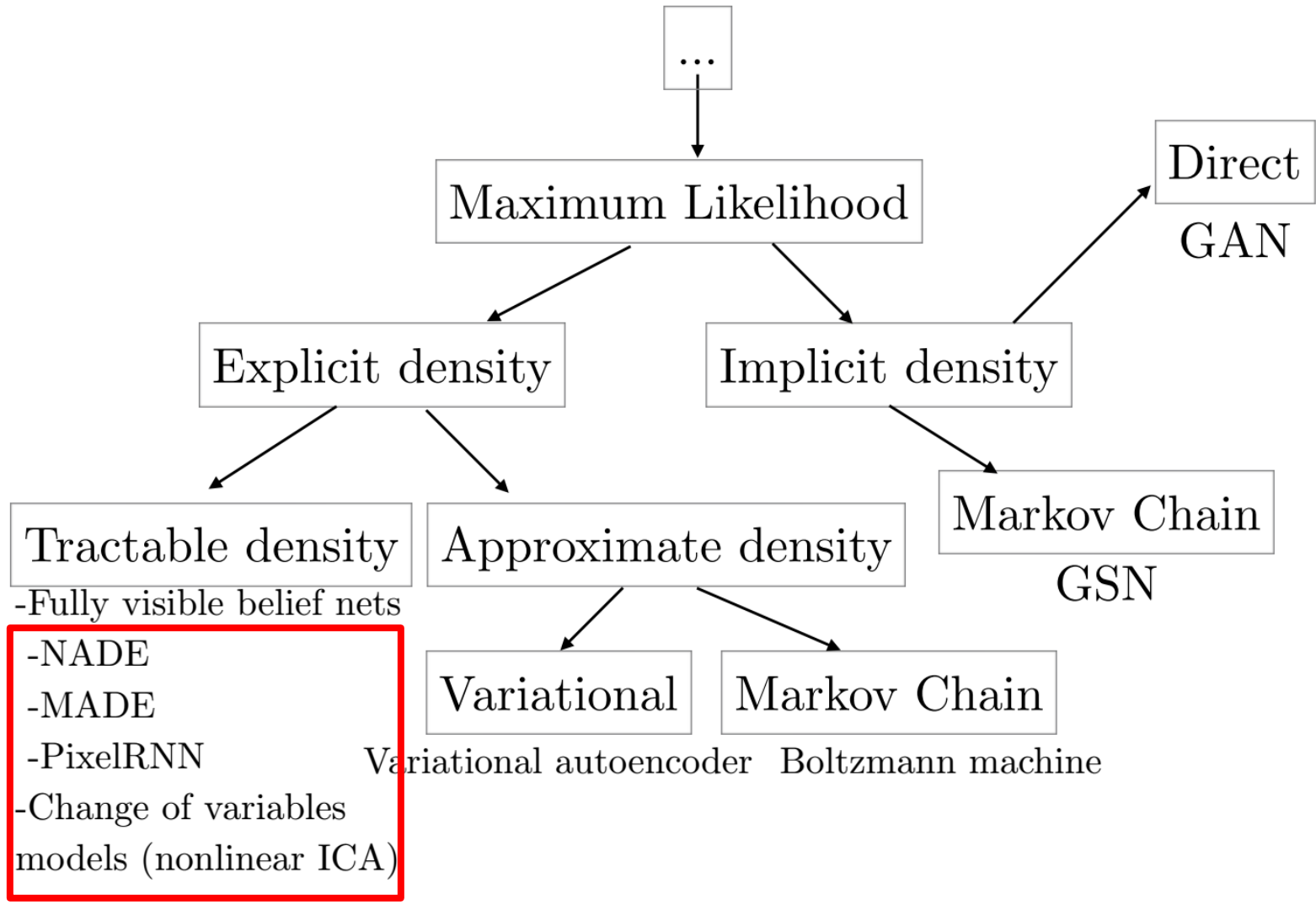
Lecture 11: Advanced Generative Models

Efstratios Gavves

Lecture overview

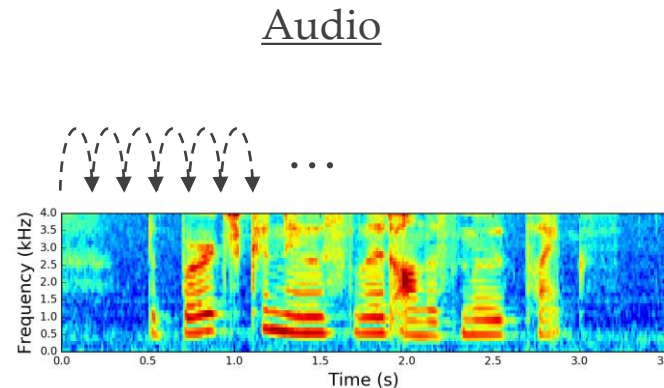
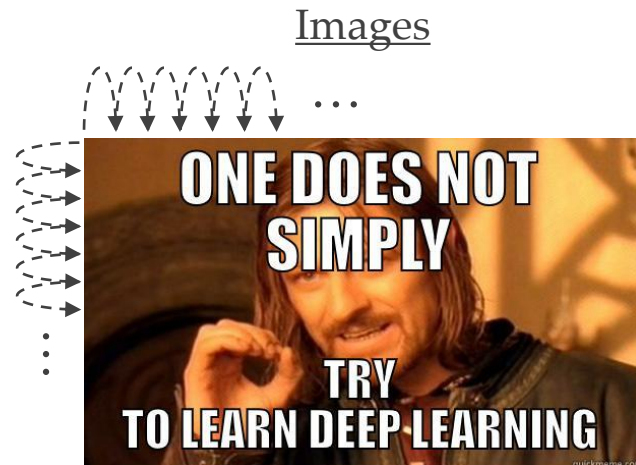
- Early autoregressive models
- Modern autoregressive models
- Normalizing flows
- Flow-based models

A map of generative models



Beyond independent dimensions

- Often, in data there is either an order or we can make up an order
 - From a generation point of view, data dimensions depend on each other



Text

...

'One does not simply try to
learn deep learning'

Decomposing likelihood of sequential data

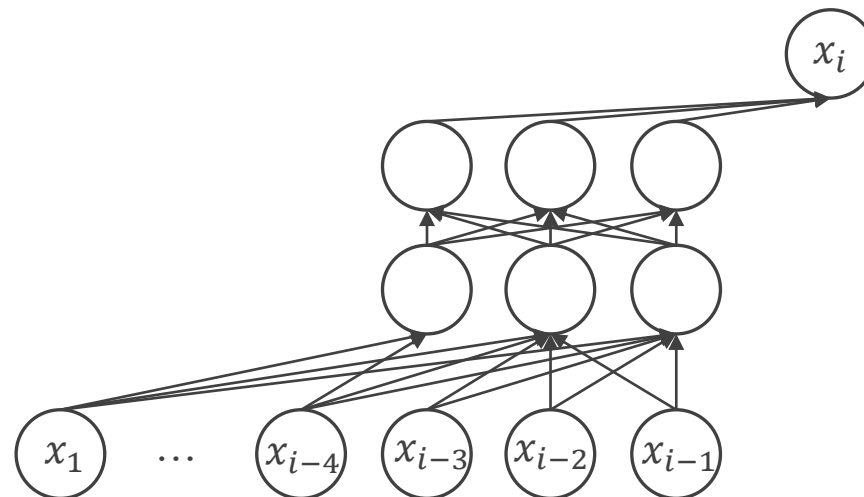
- If $\mathbf{x} = [x_1, x_2, \dots, x_d]$ is sequential, $p(\mathbf{x})$ decomposes with chain rule of probabilities

$$p(\mathbf{x}) = p(x_1) \cdot p(x_2|x_1) \cdot p(x_3|x_1, x_2) \cdot \dots \cdot p(x_d|x_1, \dots, x_{d-1}) = \prod_{i=1}^d p(x_i|x_{<i})$$

- If \mathbf{x} is *not* sequential, we can assume an artificial order
 - *e.g.*, the order with which pixels make (generate) an image
 - This can create artificial bias, however

Deep networks to model conditional likelihoods

- Model the conditional likelihoods with deep neural networks
 - Logistic regression (Frey et al., 1996), Neural nets (Bengio and Bengio, 2000)
 - *E.g.*, learn a deep net to generate one pixel at a time given past pixels
- The learning objective is to maximize the log-likelihood $\log p(\mathbf{x})$
 - If each conditional is tractable, $\log p(\mathbf{x})$ is tractable
 - Model conditional probabilities directly and with no partition functions Z



Neural Autoregressive Density Estimation

- Inspired by RBMs but with tractable density estimation
 - Each conditional modelled with sigmoidal neural net like in RBMs
- Parameter matrix W maps past inputs $\mathbf{v}_{<i}$ to hidden feature \mathbf{h}_i
- Parameter matrix V generates pixel v_i given the hidden feature \mathbf{h}_i
$$p(v_i | \mathbf{v}_{<i}) = \sigma(b_i + (V^T)_{i,\cdot} \mathbf{h}_i)$$
$$\mathbf{h}_i = \sigma(\mathbf{c} + W_{\cdot, <i} \mathbf{v}_{<i})$$
- Teacher forcing
 - During training use ground truth past inputs $\mathbf{v}_{<i}$
 - During testing use predicted past inputs $\hat{\mathbf{v}}_{<i}$

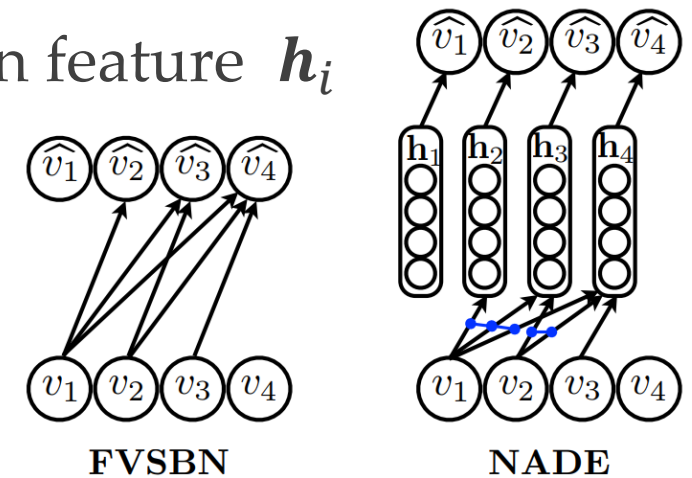


Figure 1: **(Left)** Illustration of a fully visible sigmoid belief network. **(Right)** Illustration of a neural autoregressive distribution estimator. \hat{v}_i is used as a shorthand for $p(v_i = 1 | \mathbf{v}_{<i})$. Arrows connected by a blue line correspond to connections with shared or tied parameters.

Larochelle and Murray, *Neural Autoregressive Distribution Estimation*

Masked Autoencoder for Distribution Estimation

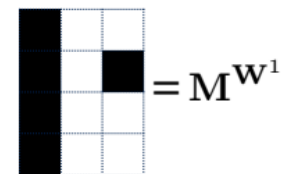
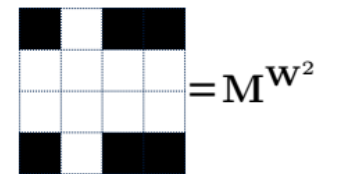
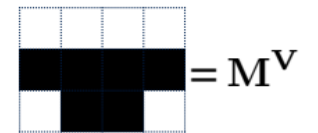
- Make an autoregressive autoencoder by setting each output x_i depend only on previous outputs $\mathbf{x}_{<i}$
 - In autoencoders the output dimensions depend on 'future' dimensions also
- Implement this by introducing a masking matrix M to multiply weights

$$h(\mathbf{x}) = g(\mathbf{b} + (W \odot M^W) \cdot \mathbf{x})$$
$$\hat{\mathbf{x}} = \sigma(\mathbf{c} + (V \odot M^V) \cdot h(\mathbf{x}))$$

For the k -th neuron the mask column is $M_{k,d} = \begin{cases} 1 & m(k) \geq d \\ 0 & \text{otherwise} \end{cases}$

And $m(k)$ is a integer between 1 and $d - 1$

Masks



Germain, Gregor, Murray, Larochelle, Masked Autoencoder for Distribution Estimation

MADE architecture

