

Lecture 11: Deep Sequential Models Efstratios Gavves

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

DEEP SEQUENTIAL MODELS - 1

• Autoregressive Models

• PixelCNN, PixelCNN++, PixelRNN

o WaveNet

o Time-Aligned DenseNets

• Let's assume we have signal modelled by an input random variable *x* • Can be an image, video, text, music, temperature measurements

• Is there an order in all these signals?

• Let's assume we have signal modelled by an input random variable *x* • Can be an image, video, text, music, temperature measurements

• Is there an order in all these signals?



• Can be an image, video, text, music, temperature measurements

o Is there an order in all these signals? Other signals and orders?



• Let's assume we have signal modelled by an input random variable *x* • Can be an image, video, text, music, temperature measurements

• Is there an order in all these signals?



o If x is sequential, there is an order: $x = [x_1, ..., x_k]$

- E.g., the order of words in a sentence
- o If x is not sequential, we can create an artificial order $x = [x_{r(1)}, ..., x_{r(k)}]$ • E.g., the order with which pixels make (generate) an image
- Then, the marginal likelihood is a product of conditionals

$$p(x) = \prod_{k=1}^{D} p(x_k | x_{j < k})$$

Different from Recurrent Neural Networks
 (a) no parameter sharing
 (b) chains are not infinite in length

• Pros: because of the product decomposition, p(x) is tractable

- Cons: because the p(x) is sequential, training is slower
- To generate every new word/frame/pixel, the previous words/frames/pixels in the order must be generated first -> no parallelism



• Sequential data is a natural fit

- Language modelling, time series, etc
- For non-sequential data they are ok, although arguably artificial
- Question: How to model the conditionals $p(x_k|x_{j < k})$
 - Logistic regression (Frey et al., 1996)
 - Neural networks (Bengio and Bengio, 2000)

Modern deep autoregressors

•NADE, MADE, PixelCNN, PixelCNN++, PixelRNN, WaveNet

NADE



Neural Autoregressive Distribution Estimation, Larochelle and Murray, AISTATS 2011

Minimizing negative log-likelihood as usual

$$\mathcal{L} = -\log p(x) = -\sum_{k=1}^{D} p(x_k | x_{< k})$$

• Then, we model the conditional as $p(x_d | x_{< d}) = \sigma(V_{d,:} \cdot h_d + b_d)$ where the latent variable h_d is defined as

$$h_d = \sigma(W_{:,d} \cdot x_d + c)$$



NADE: Training & Testing

o "Teacher forcing" training



NADE Visualizations







Binarized MNIST samples (NADE)

Binarized MNIST samples (DeepNADE)

Binarized MNIST samples (ConvNADE)

• **Question:** How could we construct an autoregressive autoencoder?

- To rephrase: How to modify an autoencoder such that each output x_k depends <u>only</u> on the previous outputs $x_{< k}$ (autoregressive property)?
 - Namely, the present k-th output \tilde{x}_k must not depend on a computational path from future inputs x_{k+1}, \ldots, x_D
 - Autoregressive: $p(x|\theta) = \prod_{k=1}^{D} p(x_k|x_{j < k}, \theta)$
 - Autoencoder: $p(\tilde{x}|x,\theta) = \prod_{k=1}^{D} p(\tilde{x}_k|x_k,\theta)$

• **Question:** How could we construct an autoregressive autoencoder?

- To rephrase: How to modify an autoencoder such that each output x_k depends <u>only</u> on the previous outputs $x_{< k}$ (autoregressive property)?
 - Namely, the present k-th output \tilde{x}_k must not depend on a computational path from future inputs x_{k+1}, \ldots, x_D
 - Autoregressive: $p(x|\theta) = \prod_{k=1}^{D} p(x_k|x_{j < k}, \theta)$
 - •Autoencoder: $p(\tilde{x}|x,\theta) = \prod_{k=1}^{D} p(\tilde{x}_k|x_k,\theta)$

• Answer: Masked convolutions!

$$h(x) = g(b + (W \odot M^{W}) \cdot x)$$

$$\tilde{x} = \sigma(c + (V \odot M^{V}) \cdot h(x))$$

Masks



Masked Autoencoder for Distribution Estimation, Germain, Mathieu et al., ICML 2015

MADE



Masked Autoencoder for Distribution Estimation, Germain, Mathieu et al., ICML 2015

UVA DEEP LEARNING COURSE – EFSTRATIOS GAVVES

PixelRNN

• Unsupervised learning: learn how to model p(x)

• Decompose the marginal

$$p(x) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$$
TO LEARN DEEP LEARNING

2

 \circ Assume row-wise pixel by pixel generation and sequential colors $R \rightarrow G \rightarrow B$ • Each color conditioned on all colors from previous pixels and specific colors in the same pix $p(x_{i,R}|x_{\leq i}) \cdot p(x_{i,G}|x_{\leq i}, x_{i,R}) \cdot p(x_{i,B}|x_{\leq i}, x_{i,R}, x_{i,G})$

• Final output is 256-way softmax

Pixel Recurrent Neural Networks, van den Oord, Kalchbrenner and Kavukcuoglu, arXiv 2016

S MP V

TRY

• How to model the conditionals? $p(x_{i,R}|x_{<i}), p(x_{i,G}|x_{<i}, x_{i,R}), p(x_{i,B}|x_{<i}, x_{i,R}, x_{i,G})$

LSTM variants
 12 layers

Row LSTMDiagonal Bi-LSTM



- Hidden state (i, j) =
 Hidden state (i-1, j-1) +
 Hidden state (i-1, j) +
 Hidden state (i-1, j+1) +
 p(i, j)
- By recursion the hidden state captures a fairly triangular region



• How to capture the whole previous context

o Pixel (i, j) =
 Pixel (i, j-1) +
 Pixel (i-1, j)

• Processing goes on diagonally

• Receptive layer encompasses entire region

Diagonal Bi-LSTM

• Propagate signal faster

• Speed up convergence



• Pros: good modelling of $p(x) \rightarrow$ nice image generation

• Half pro: Residual connections speeds up convergence

• Cons: still slow training, slow generation





Figure 1. Image completions sampled from a PixelRNN.

• Unfortunately, PixelRNN is too slow

Solution: replace recurrent connections with convolutions
 Multiple convolutional layers to preserve spatial resolution

• Training is much faster because all true pixels are known in advance, so we can parallelize

 \circ Generation still sequential (pixels must be generated) \rightarrow still slow



PixelCNN

Stack of masked



Pixel Recurrent Neural Networks, van den Oord, Kalchbrenner and Kavukcuoglu, arXiv 2016

• Use masked convolutions again to enforce autoregressive relationships





PixelCNN – Pros and Cons

Cons: Performance is worse than PixelRNN Why?

o Cons: Performance is worse than PixelRNN

○Why?

• Not all past context is taken into account

• New problem: the cascaded convolutions create a "blind spot"

o Because of

- (a) the limited receptive field of convolutions and
- (b) computing all features at once (not sequentially)
- \rightarrow cascading convolutions makes current pixel not depend on <u>all</u> previous
- \rightarrow blind spot



Use two layers of convolutions stacks

 Horizontal stack: conditions on current row and takes as input the previous layer output and the vertical stack

• Vertical stack: conditions on all rows above current pixels

• Also replace ReLU with a $tanh(W * x) \cdot \sigma(U * x)$



PixelCNN - Generations

o Coral reef



o Sorrel horse



o Sandbar



o Lhasa Apso



o Improving the PixelCNN model

- Replace the softmax output with a discretized logistic mixture lihelihood
- Softmax is too memory consuming and gives sparse gradients
- Instead, assume logistic distribution of intensity and round off to 8-bits
- Condition on whole pixels, not pixel colors
- Downsample with stride-2 convs to compute long-range dependencies
- Use shortcut connections
- o Dropout
- $^{\circ}$ PixelCNN is too powerful a framework ightarrow can onverfit easily

PixelCNN++: Improving the PixelCNN with Discretized Logistic, Salimans, Karpathy, Chen, Kingma, ICLR 2017



PixelCNN++ - Generations



• A standard VAE with a PixelCNN generator/decoder

 O Be careful. Often the generator is so powerful, that the encoder/inference network is ignored ← Whatever the latent code z there will be a nice image generated

PixelVAE: A Latent Variable Model for Natural Images, Gulrajani et al., ICLR 2017



PixelVAE - Generations



64x64 LSUN Bedrooms

64x64 ImageNet

PixelVAE - Generations

Varying top latents

Varying bottom latents

Varying pixel-level noise





WaveNet: A Generative Model for Raw Audio, van den Oord et al., arXiv 2017

o Inspired by PixelRNN and PixelCNN

• Fully convolutional neural network

• Use dilated convolutions

o <u>Samples</u>



WaveNet architecture



• Video Time: Properties, Encoders and Evaluation • A. Ghodrati, E. Gavves, C. Snoek, BMVC 2018

• How to model image sequences optimally?



Video Time: Properties, Encoders and Evaluation, Ghodrati, Gavves, Snoek, BMVC 2018

Properties of video

A standard VAE with a PixelCNN



Forward



Backward

Video properties

Natural order (+)



Reverse order (-)



• Temporal continuity

o Temporal asymmetry



"Pretending to put something into something"

• Temporal causality



LSTMs learn transitions between subsequent states

3D convolutions learn spatiotemporal patterns within a video





Autoregressive-like: each new frame directly depends on all previous frames
 They are not generative

 \circ No parameter sharing \rightarrow like ConvNets, unlike RNNs

o Sequential → like RNNs, unlike ConvNets



Better latent space



Summary

Autoregressive Models
PixelCNN, PixelCNN++, PixelRNN
WaveNet
Time-Aligned DenseNets

UVA DEEP LEARNING COURSE EFSTRATIOS GAVVES DEEP SEQUENTIAL MODELS - 52