Deep Generative Models

Chapter 3: Generation by Explicit Distribution Learning

Ali Bereyhi

ali.bereyhi@utoronto.ca

Department of Electrical and Computer Engineering University of Toronto

Summer 2025

1/26

Overview and LMs

We have already worked a lot with AR models: in Chapter 1

we built various LMs *vor* they are AR models

People have tried AR modeling for other modalities as well

- It has been largely developed for visual modalities like images
- Various challenges like causality are faced there too
 - └→ Concepts like masked decoding or positional encoding were developed there
 - → Modern LMs indeed borrow lots of those notions

In this section, we go through some major computational AR models

RNN-based AR Models

Recurrent AR models are one of primary computational AR models: we have already seen the recurrent LMs

In these models, we use RNN to build conditional distribution

- **1** The model samples x_1 from a learned prior
- **2** Given x_i and the state of RNN at i 1, the model
 - $\, {\scriptstyle {\scriptstyle {\rm L}}} \,$ updates its state to represent the context at i
 - \downarrow computes Softmax-activated output of context to estimate $P(x_{i+1}|x_{< i+1})$
- 3 The model iterates till d entries are over

Key Observation

We cannot do parallel processing vor training is also slow











Recurrent AR Models

PixelRNN: Visual Recurrent AR Model







Recurrent AR Models

PixelRNN: Visual Recurrent AR Model



Recurrent AR Models

PixelRNN: Visual Recurrent AR Model







PixelRNN: Training

We can readily train this model via MLE



PixelRNN: Generation

It's good to think how we can generate with this model



Recurrent AR Models

PixelRNN: Multi-channel

For multi-channel data, we typically first process in depth

→ We expect that RGB pixels are heavily correlated



Channel

Channel 2

Channel 3

PixelRNN: Row Processing



Basic PixelRNN goes through images row-wise \rightsquigarrow processing time $\propto d^2$

- For generation we need to complete these d^2 steps
- For training we can treat each row independently pprox truncated BPTT
 - → Faster at the price of capturing less spatial correlation

PixelRNN: Diagonal Processing



Alternatively, we can process each pixel conditioned on its top and left pixels

- We can process multiple pixels in parallel
- Total processing time for a single image is 2d-1 or ∞d
 - → Fast with less loss in capturing spatial correlation

PixelRNN: Diagonal Processing - Reduced AR Modeling

x_1	x_2		
x_7			

Attention

Note that we get the benefit of parallel processing due to reduced AR modeling

$$P(x_2|x_1) P(x_7|x_1, x_2) \approx P(x_7|x_1) P(x_7|x_1)$$

Deep Generative Models

Convolutional AR Models

When we work with visual data: it sounds natural that we work with CNNs

- + But, CNNs cannot capture context? We need some sort of state!
- But, we could extract context with Transformers
 - They only consist of bunch of MLP-like layers
 - We made context without an explicit state

This is in fact a well-developed idea!

Spatial Context via Convolution

We can compute a distribution for a single pixel using multiple kernels



Valid Conditional Distribution

- + Wait a moment! The conditional distribution however depends on all pixels in the receptive field! Including the pixel itself! Right?!
- Yes! You are right!
- + Then, this cannot be a valid conditional distribution!
 - We need to compute distribution conditioned to what we had before!
 - This distribution depends on the pixel itself and some future ones!
- We can resolve this using masked filter
 - We did the same in Transformers!

Masking Kernels

We mask kernels at the pixel position and the future positions





























 _	-		

 $\times 256$





Deep Generative Models















		$\times 256$	





Short Range Conditioning

Short Range Impact of Masked Convolution

The computed conditional distributions are valid; however, they are conditioned only on nearby pixels within the receptive field of kernel

$x_{i-1,\ell-1}$	$x_{i-1,\ell}$	$x_{i-1,\ell+1}$	
$x_{i,\ell-1}$	$x_{i,\ell}$		

In terms of AR modeling, we assume

$$P(x_{i,\ell}|x_{< i,\ell}) \approx P(x_{i,\ell}|x_{\text{neighbors}})$$

which weakens spatial correlation capturing

Recap: Receptive Field of Deep CNN

Recap: Deep Convolution

Recall as we go deeper in CNN vvv effective receptive field expands



Solution: Deep Masked Convolution

Using deep masked CNN vor effective receptive field covers all previous pixels

Deep Masked Convolution



As we go deep with masked convolutional layers

- receptive field expands to the whole image scene
- since all kernels are masked, the final receptive field is also masked
- this extracts an autoregressive context that we need

Deep Masked Convolution



As we go deep with masked convolutional layers

- receptive field expands to the whole image scene
- since all kernels are masked, the final receptive field is also masked
- this extracts an autoregressive context that we need

Deep Masked Convolution

As we go deep with masked convolutional layers

- receptive field expands to the whole image scene
- since all kernels are masked, the final receptive field is also masked
- this extracts an autoregressive context that we need

PixelCNN: AR Model with Deep Masked Convolution



PixelCNN: Few Notes



In practical implementation of PixelCNN

- Typically the network is quite deep > 15
 - └→ To guarantee accurate conditioning
- Deep networks can quickly experience vanishing gradients
 - → We use skip connextions and gating to reduce vanishing gradients
- Final layer extracts 256 features for each pixel
 - Like PixelRNN, we typically condition RGB channels right after each other

PixelCNN: Training

We can train PixelCNN via MLE

 $\begin{array}{l} \underline{\operatorname{Train_PixelCNN}(\mathbb{D}: \operatorname{dataset}):} \\ 1: \ \text{for multiple epochs do} \\ 2: \ Sample \ a \ batch \ of \ n \ images \ \left\{x^j: j=1,\ldots,n\right\} \\ 3: \ \ \text{for } j=1:n \ \text{do} \\ 4: \ \ \mathbf{p}_{1:d}^j \leftarrow \operatorname{MaskedCNN}(x^j) \\ 5: \ \ \text{end for} \\ 6: \ \ Update \ model \ using \ \operatorname{Opt_avg}\left\{\sum_i \nabla \log \mathbf{p}_i^j \left[x_i^j\right]\right\} \\ 7: \ \text{end for} \\ 8: \ \ \text{return \ trained \ model} \end{array}$

Parallel Processing

Using masked convolution, we can process all pixels in parallel while we train

PixelCNN: Generation

Generation is still autoregressive

```
\begin{array}{l} \underline{\text{Sample_PixelCNN():}}\\ \hline \textbf{1: Set } x^{(0)} = \varnothing\\ 2: \text{ for } i = 1, \dots, d \text{ do}\\ 3: \quad \textbf{p}_{1:d} \leftarrow \text{MaskedCNN}(x^{(i-1)})\\ 4: \quad x_i \leftarrow \text{multinomial}(\textbf{p}_i, \ \texttt{#smpl=1})\\ 5: \quad \text{Update } x^{(i)} \leftarrow \ \text{Replace } x_i \ \text{in } x^{(i-1)}\\ 6: \ \text{end for}\\ 7: \ \text{return } x^{(d)} \end{array}
```

#prior sampling

Transformer-Based AR

We can use the same approach as in LMs to extract context via transformers



This is exactly LM with time \equiv position and tokens \equiv pixels

Deep Generative Models

ImageGPT: Extending LMs to Generate Pixels

A good example is ImageGPT:

it interprets pixels as tokens \rightsquigarrow image is a long sequence of tokens

- Vocalbulary has only 256 tokens
 - → We can patch multiple pixels together to make larger vocabulary
 - → With patch-pixels, we can use tokenization to reduce vocabulary size
- Each pixel-patch is embedded with an embedding layer
 - → This embedding can be in general smaller than vocabulary size
 - → It will give us a numerical representation of the image patches
- As we use self-attention, the model learns spatial correlation by itself
 - → This is more robust as compared to deep masked convolution
 - → It however needs more data to start capturing spatial pattern

Take a look at ImageGPT Paper

Wrap Up and Final Notes

AR modeling gives an efficient way to computationally model data distribution

→ It deploys autoregressive sampling leading from slow generation

We can realize them by

- Recurrent networks, e.g., LSTM
 - → It explicitly conditions distribution through model state
 - → We do not benefit from parallel processing
- Masked networks, e.g., masked CNN or Transformers
 - → It conditions causally through masking
 - → We enjoy the luxury of parallel processing

AR modeling are widely considered slow of visual modalities

→ They are though getting revisited due to some recent advancement

Next Stop: Energy-based Models!

26/26