Predictive Learning for Video Prediction

Mingsheng Long

Tsinghua University



Today's Al

Today, AI = Supervised (Deep) Learning





"Real" Al

But Supervised Learning is Insufficient for "Real" Al

- Most of animals and humans learning is unsupervised, through interaction with the world
- We learn how the world works by observing it
 - We learn many simple things: depth and 3dimensionality, gravity, object permanence, ...
- We build models of the world through predictive unsupervised learning
- World models give us "common sense"









Predictive Learning



How Much Information does the Machine Need to Predict?

Pure Reinforcement Learning (cherry)

- The machine predicts a scalar reward given once in a while.
- ► A few bits for some samples

Supervised Learning (icing)

- The machine predicts a category or a few numbers for each input
- ▶ 10→10,000 bits per sample

Unsupervised/Predictive Learning (génoise)

- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos
- Millions of bits per sample



Unsupervised Learning is the Dark Matter (or Dark Energy) of AI



Video Prediction: Learning Physics

Learning Physics (PhysNet) [Lerer, Gross, Fergus, ICML'16]

- ConvNet predicts the trajectories of falling blocks
- Uses the Unreal game engine hooked up to Torch.





Video Prediction: Weather Nowcasting

March 4th, Jiangxi Province, China





Predictive Learning

Predictive Learning



Technical Challenges





CNN-based methods

- 3D CNNs [Vondrick 2016]

RNN-based methods

- ConvLSTM network [Shi 2015]

- Video Pixel Networks [Kalchbrenner 2017]

- PredRNN / PredRNN++ [Wang 2017, Wang 2018]



VideoGAN



Generating Videos with Scene Dynamics. Vondrick et al. NIPS 2016.



Warm-Up: Seq-to-Seq LSTMs



[Srivastava et al. Unsupervised Learning of Video Representations using LSTMs. ICML 2015]



Warm-Up: Seq-to-Seq LSTMs





Long Short-Term Memory (LSTM)



LSTM

Input, forget, output gates: $i_{t} = \sigma(W_{xi}x_{t} + W_{hi}h_{t-1} + b_{i})$ $f_{t} = \sigma(W_{xf}x_{t} + W_{hf}h_{t-1} + b_{f})$ $o_{t} = \sigma(W_{xo}x_{t} + W_{ho}h_{t-1} + b_{o})$

Cell state: $g_t = \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c)$ $c_t = f_t \odot c_{t-1} + i_t \odot g_t$

I. All variables are ID vectors2. Dimensions are Permutable (Spatial Information Loss)

Hidden layer: $h_t = o_t \odot \tanh(c_t)$



Convolutional LSTM (ConvLSTM)



[Convolutional LSTM Network. Shi et al. NIPS 2015.]



ConvLSTM





State-to-state Convolution: future states \rightarrow larger receptive field

[Convolutional LSTM Network. Shi et al. NIPS 2015.]



ConvLSTM Network



[Convolutional LSTM Network. Shi et al. NIPS 2015.]



PixelRNN



• Then maximize likelihood of training data $\{x\}$

A Oord et al. Pixel Recurrent Neural Networks. ICML 2016.



PixelRNN



A Oord et al. Pixel Recurrent Neural Networks. ICML 2016.



Video Pixel Networks (VPN)



Video Pixel Networks. Kalchbrenner et al. ICML 2017.



PredRNN

Zigzag Info Flow

- Deeper transitions for short-term (fast flow) dependencies
- Larger receptive field across adjacent states for sudden change



PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs. Yunbo Wang et al. NIPS 2017.



PredRNN

Zigzag Info Flow

- Zigzag flow makes network deeper in time \rightarrow gradient vanishing
- Maintain temporal memory $(C_t) \rightarrow$ Preserve long-term gradients



PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs. Yunbo Wang et al. NIPS 2017.



PredRNN

Dual memory

- Spatiotemporal Memory (in yellow): an external memory (M_t)
- Temporal Memory (in black): the conventional memory cell (C_t^{\prime})



PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs. Yunbo Wang et al. NIPS 2017.



PredRNN++

Zigzag flow: more straightforward supervisions to the memory



(a) Stacked LSTMs (b) Deep Transition LSTMs

But, for the short term -> deeper-in-time networks -> gradient vanishing -> bad long-term modeling capability



PredRNN++: Longer path for short-term dynamics



[PredRNN++. Yunbo Wang et al. ICML 2018.]



PredRNN++: Longer path for short-term dynamics



A dual memory structure, C^k_t: the conventional LSTM memory
 M^k_t: a zigzag delivered memory (vertically and horizontally)
 C^k_t → M^k_t: add more non-linear layers to recurrent transitions from one time step to the next.

$$\begin{pmatrix} g_t \\ i_t \\ f_t \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \end{pmatrix} W_1 * [\mathcal{X}_t, \mathcal{H}_{t-1}^k, \mathcal{C}_{t-1}^k]$$

$$\mathcal{C}_t^k = f_t \odot \mathcal{C}_{t-1}^k + i_t \odot g_t$$

$$\begin{pmatrix} g_t' \\ i_t' \\ f_t' \end{pmatrix} = \begin{pmatrix} \tanh \\ \sigma \\ \sigma \end{pmatrix} W_2 * [\mathcal{X}_t, \mathcal{C}_t^k, \mathcal{M}_t^{k-1}]$$

$$\mathcal{M}_t^k = f_t' \odot \tanh (W_3 * \mathcal{M}_t^{k-1}) + i_t' \odot g_t'$$

$$o_t = \tanh (W_4 * [\mathcal{X}_t, \mathcal{C}_t^k, \mathcal{M}_t^k])$$

$$\mathcal{H}_t^k = o_t \odot \tanh (W_5 * [\mathcal{C}_t^k, \mathcal{M}_t^k])$$

[PredRNN++. Yunbo Wang et al. ICML 2018.]



PredRNN++: Shorter path for long-term relations



[PredRNN++. Yunbo Wang et al. ICML 2018.]



Moving MNIST: Frequent Occlusions

		MNI	MNIST-3				
Model	10 tim	E STEPS	30 tim	E STEPS	10 TIME STEPS		
	SSIM	MSE	SSIM	MSE	SSIM	MSE	
FC-LSTM (SRIVASTAVA ET AL., 2015A)	0.690	118.3	0.583	180.1	0.651	162.4	
CONVLSTM (SHI ET AL., 2015)	0.707	103.3	0.597	156.2	0.673	142.1	
TrajGRU (Shi et al., 2017)	0.713	106.9	0.588	163.0	0.682	134.0	
CDNA (FINN ET AL., 2016)	0.721	97.4	0.609	142.3	0.669	138.2	
DFN (DE BRABANDERE ET AL., 2016)	0.726	89.0	0.601	149.5	0.679	140.5	
VPN* (KALCHBRENNER ET AL., 2017)	0.870	64.1	0.620	129.6	0.734	112.3	
PREDRNN (WANG ET AL., 2017)	0.867	56.8	0.645	112.2	0.782	93.4	
CAUSAL LSTM	0.882	52.5	0.685	100.7	0.795	89.2	
CAUSAL LSTM (VARIANT: SPATIAL-TO-TEMPORAL)	0.875	54.0	0.672	103.6	0.784	91.8	
PredRNN + GHU	0.886	50.7	0.713	98.4	0.790	88.9	
CAUSAL LSTM + GHU (FINAL)	0.898	46.5	0.733	91.1	0.814	81.7	

Inputs	Ð	P	Ð	80	80	80	e B B	060	Oyo	Ŷo	Inputs	Nm	Nm	23	ry J	z	z	Ð	3	R.	Negy
Targets	G	ß	Ø	8	Oan	00	60	60	08	08	Targets	Nn	Nn	Nm	23	2 M	5° S	N M	N M	Nu	N
PredRNN ++	G	Ð	Ø	B	B	00	000	00	08	08	PredRNN ++	N	Nm	Nm	2 3	N 33	N 3	N 33	Nw	Nw	N
PredRNN	Ð	Ŷ	8	ð	Ord	0.50	50	08	05	06	PredRNN	N	M	NM	Ng	23	<i>v</i> 3	2 S S	23	23	23
VPN baseline	Ŷ	Ð	Ð	B	50	30	^ъ	° 0	03	03	VPN baseline	Nay	Nam	Nem	الله على	N 37)	(m 11	N 33	n m	NM	Nm
ConvLSTM	Po	Ð	ŝ	Ð	0.0	0.0	0	05	05	05	ConvLSTM	Copo.	nd m	NIT	Nin	(n) P	(u) PV	n m	NM	n M	23
TrajGRU	P	Ê	B	ð	ð	00	000	4	4	4	TrajGRU	Non	Nen	N .m	n m	Un N	ы. N	5° 39	N	5	2



Moving MNIST: Frequent Occlusions







Predictive Learning

Real Videos: Predictable Movement

Model	PSNR	SSIM
Ours	28.47	0.865
PredRNN [Wang et al., 2017]	27.55	0.839
MCnet [Villegas et al., 2017]	25.95	0.804
ConvLSTM [Shi et al., 2015]	23.58	0.712
DFN [De Brabandere et al., 2016]	27.26	0.794

Inputs				Targe	ets a	nd pı	redic	tions	5	Inputs Targets and predictions									
t=3	t=6	t=9	t=12	t=15	t=18	t=21	t=24	t=27	t=30	t=3	t=6	t=9	t=12	t=15	t=18	t=21	t=24	t=27	t=30
+	X	X	k	1	8	X	K	k		1	k	X	1	k	X	1	1	t	1
Pre	dRNN	+	k	1	8	λ	X	k	!	Pre	dRNN+	-+	1	1	1	1	1	1	1
Pre	dRNN		k	1	A	X	K			PredRNN			1	*	1	1	+	1	1
MCn	et		k	1	8	X	A	!		MCnet			1	Å	X	1	8	4	
Con	VLSTN	1	k	1	k	X	A	1	1	ConvLSTM			1	Å	x	1	A	A	1
DFN			k	1	k	K	٨	1		DFN			1	1	1	4	1	1	1



Weather Forecasting: A Case Study





 x_{t-2}

 x_{t-1}

 x_t

I. Layer Normalization -> Convergence

After convolutions but right before nonlinearity

- Cutting memory usage **87.5%** by saving mini-batch
- Cutting training time 75%
- Increase accuracy 20%

Layer Normalization

Recap: batch normalization

$$\bar{a}_{i}^{l} = \frac{g_{i}^{l}}{\sigma_{i}^{l}} (a_{i}^{l} - \mu_{i}^{l}) \qquad \mu_{i}^{l} = E_{x \sim P(x)} [a_{i}^{l}] \qquad \sigma_{i}^{l} = \sqrt{E_{x \sim P(x)} [(a_{i}^{l} - \mu_{i}^{l})^{2}]}$$

Layer normalization in RNNs

$$\boldsymbol{h}^{t} = f[\frac{\boldsymbol{g}}{\sigma^{t}} \odot(\boldsymbol{a}^{t} - \boldsymbol{\mu}^{t}) + \boldsymbol{b}] \qquad \boldsymbol{\mu}^{t} = \frac{1}{C} \sum_{i=1}^{C} a_{i}^{t} \qquad \boldsymbol{\sigma}^{t} = \sqrt{\frac{1}{C} \sum_{i=1}^{C} (a_{i}^{t} - \boldsymbol{\mu}^{t})^{2}}$$







II. Transpose Convolution -> OOM

Many intermediate activations -> memory explosion

Resolution 400*400 -> Memory usage 18.31GB

Transpose convolution -> 7.15GB





Predictive Learning

II. Sampling -> OOM





II. Multiple GPUs -> OOM





III. Loss Function -> Blurry





IV. Data Augmentation -> Overfitting



- **IO** times training set
- **IO%** accuracy boost





V. Scheduled Sampling



Bengio, S. et al. "Scheduled Sampling for Sequence Prediction with Recurrent Neural Networks." NIPS 2015.



Weather Forecasting

Algorithm	Inventor	Year	MSE/frame
ConvLSTM	СИНК	2015	10.48
CDNA	OpenAl	2016	9.11
VPN	DeepMind	2017	8.49
TrajGRU	СИНК	2017	8.43
PredRNN	Tsinghua	2017	7.54
PredRNN++	Tsinghua	2018	6.38





Weather Forecasting



Precipitation Nowcasting

Typhoon Monitoring



Thank You Questions?

Mingsheng Long <u>mingsheng@tsinghua.edu.cn</u> <u>http://ise.thss.tsinghua.edu.cn/~mlong</u>