

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532L): **Multimodal Learning with Vision, Language and Sound**

Lecture 8: RNNs (cont.) + Applications



Course Logistics

- Assignment 1 grades are (finally) available
 - Part 1 solutions are available on Piazza
 - Part 2 solutions will be out soon
- Assignment 3 is out, due Wednsday, February 7th
- Questions on Assignment 3?

- **Projects** pitches we will do during next class

- Paper choices (google form) will is due Wednsday THIS week (Jan 31st)

RNNS: Review Key Enablers:

- Parameter sharing in computational graphs
- "Unrolling" in computational graphs
- Allows modeling arbitrary length sequences!





RNNS: Review Key Enablers:

- Parameter sharing in computational graphs
- "Unrolling" in computational graphs
- Allows modeling arbitrary length sequences!







RNNS: Review Key Enablers:

- Parameter sharing in computational graphs
- "Unrolling" in computational graphs
- Allows modeling arbitrary length sequences!







RNNs: Review **Key Enablers:**

- Parameter sharing in computational graphs
- "Unrolling" in computational graphs
- Allows modeling arbitrary length sequences!







Vanishing Or Exploding Gradients

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$
$$c_t = f \odot c_{t-1} + i \odot g$$
$$h_t = o \odot \tanh(c_t)$$



Uninterrupted gradient flow!



RNNs: Review **Key Enablers:**

- Parameter sharing in computational graphs
- "Unrolling" in computational graphs
- Allows modeling arbitrary length sequences!
- or Squared Loss (regression)





Loss functions: often cross-entropy (for classification); could be max-margin (like in SVM)





Long-Short Term Memory (LSTM)



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)



* slide from Dhruv Batra

Long-Short Term Memory (LSTM)

Cell state / **memory**



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

* slide from Dhruv Batra

LSTM Intuition: Forget Gate

Should we continue to **remember** this "bit" of information or not?



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$f_t = \sigma \left(W_f \cdot [h_{t-1}, x_t] + b_f \right)$

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Input Gate

Should we **update** this "bit" of information or not? If yes, then what should we **remember**?



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$i_t = \sigma \left(W_i \cdot [h_{t-1}, x_t] + b_i \right)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Memory Update



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Forget what needs to be forgotten + memorize what needs to be remembered

$C_t = f_t * C_{t-1} + i_t * C_t$

* slide from Dhruv Batra



LSTM Intuition: Output Gate

Should we output this bit of information (e.g., to "deeper" LSTM layers)?



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$o_t = \sigma \left(W_o \left[h_{t-1}, x_t \right] + b_o \right)$ $h_t = o_t * \tanh(C_t)$

* slide from Dhruv Batra

, Dotr

LSTM Intuition: Additive Updates



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Uninterrupted gradient flow!

* slide from Dhruv Batra



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

Uninterrupted gradient flow!

* slide from Dhruv Batra



LSTM Variants: with Peephole Connections

Lets gates see the cell state / memory



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$f_t = \sigma \left(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f \right)$$

$$i_t = \sigma \left(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i \right)$$

$$o_t = \sigma \left(W_o \cdot [C_t, h_{t-1}, x_t] + b_o \right)$$

* slide from Dhruv Batra

, Dotr

LSTM Variants: with Peephole Connections

Lets gates see the cell state / memory



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$f_{t} = \sigma \left(W_{f} \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_{t}] + b_{f} \right)$$

$$i_{t} = \sigma \left(W_{i} \cdot [\boldsymbol{C_{t-1}}, h_{t-1}, x_{t}] + b_{i} \right)$$

$$o_{t} = \sigma \left(W_{o} \cdot [\boldsymbol{C_{t}}, h_{t-1}, x_{t}] + b_{o} \right)$$

* slide from Dhruv Batra

, Dotr

LSTM Variants: with Coupled Gates

Only memorize new information when you're forgetting old



Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$C_t = f_t * C_{t-1} + (1 - f_t) * \tilde{C}_t$

* slide from Dhruv Batra

Gated Recurrent Unit (GRU)

No explicit memory; memory = hidden output



z = memorize new and forget old

Image Credit: Christopher Olah (http://colah.github.io/posts/2015-08-Understanding-LSTMs/)

$$z_t = \sigma \left(W_z \cdot [h_{t-1}, x_t] \right)$$
$$r_t = \sigma \left(W_r \cdot [h_{t-1}, x_t] \right)$$
$$\tilde{h}_t = \tanh \left(W \cdot [r_t * h_{t-1}, x_t] \right)$$
$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

* slide from Dhruv Batra

LSTM/RNN Challenges

- LSTM can remember some history, but not too long - LSTM assumes data is regularly sampled



Phased LSTM

Gates are controlled by **phased** (periodic) **oscillations**



[Neil et al., 2016]



Bi-directional RNNs/LSTMs

$$y_t = W_{hy}h_t + b_y$$

$$h_t = f_W(h_{t-1}, x_t)$$

 $h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + b_h)$

Bi-directional RNNs/LSTMs

$$y_t = W_{hy} [\overrightarrow{h}_t, \overleftarrow{h}_t]^T + b$$

$$\overrightarrow{h}_{t} = f_{\overrightarrow{W}}(\overrightarrow{h}_{t-1}, x_{t})$$

$$\overleftarrow{h}_{t} = f_{\overleftarrow{W}}(\overleftarrow{h}_{t+1}, x_{t})$$

$$\overrightarrow{h}_{t} = \tanh(\overrightarrow{W}_{hh}\overrightarrow{h}_{t-1} + \overrightarrow{W}_{xh}x_{t} + \overrightarrow{h}_{t} = \tanh(\overleftarrow{W}_{hh}\overleftarrow{h}_{t+1} + \overleftarrow{W}_{xh}x_{t} + \overrightarrow{h}_{t})$$

 b_y

 $\overrightarrow{b}_{h})$ $\overleftarrow{b}_{h})$



Consider a translation task: This is one of the first neural translation models



French Decoder English Encoder

English Encoder

Consider a translation task with a bi-directional encoder of the source language







Consider a translation task with a bi-directional encoder of the source language





Build a small neural network (one layer) with softmax output that takes (1) everything decoded so far and (encoded by previous decoder state Zi) (2) encoding of the current word (encoded by the hidden state of encoder hj) and predicts relevance of every source word towards next translation

Consider a translation task with a bi-directional encoder of the source language

Cho et al., 2015]







Build a small neural network (one layer) with softmax output that takes (1) everything decoded so far and (encoded by previous decoder state Zi) (2) encoding of the current word (encoded by the hidden state of encoder hj) and predicts relevance of every source word towards next translation

Consider a translation task with a bi-directional encoder of the source language

https://devblogs.nvidia.com/introduction-neural-machine-translation-gpus-part-3/



Cho et al., 2015]





[Cho et al., 2015]



Regularization in RNNs

Standard dropout in recurrent layers of **long term memory**!

Standard dropout in recurrent layers does not work because it causes loss of

Regularization in RNNs

long term memory!

- Dropout in input-to-hidden or hidden-to-output layers [Zaremba et al., 2014]
- Apply dropout at sequence level (same zeroed units for the entire sequence)
- Dropout only at the cell update (for LSTM and GRU units) [Semeniuta et al., 2016]
- Enforcing norm of the hidden state to be similar along time [Krueger & Memisevic, 2016]
- Zoneout some hidden units (copy their state to the next tilmestep) [Krueger et al., 2016]

Standard dropout in recurrent layers does not work because it causes loss of















Training Objective: Predict the next word (cross entropy loss)



Testing: Sample the full sequence





Training Objective: Predict the next word



Testing: Sample the full sequence



Slowly move from Teacher Forcing to Sampling







Approach vs Metric

Baseline **Baseline with Dropout Always Sampling** Scheduled Sampling **Uniform Scheduled Sampling** Baseline ensemble of 10 Scheduled Sampling ensemble of

Baseline: Google NIC captioning model

Baseline with Dropout: Regularized RNN version

Always sampling: Use sampling from the beginning of training

Scheduled sampling: Sampling with inverse Sigmoid decay

Uniformed scheduled sampling: Scheduled sampling but uniformly

Microsoft COCO developement set			
	BLEU-4	METEOR	CIDER
	28.8	24.2	89.5
	28.1	23.9	87.0
	11.2	15.7	49.7
	30.6	24.3	92.1
5	29.2	24.2	90.9
	30.7	25.1	95.7
of 5	32.3	25.4	98.7

Sequence Level Training

- During training objective is different than at test time
- **Training:** generate next word given the previous
- **Test:** generate the entire sequence given an initial state

Optimize directly evaluation metric (e.g. BLUE score for sentence generation)

Set the problem as a Reinforcement Learning:

- RNN is an Agent
- Policy defined by the learned parameters
- Action is the selection of the next word based on the policy Reward is the evaluation metric

[Ranzato et al., 2016]
Let us look at some actual practical uses of RNNs

Applications: Skip-thought Vectors



word2vec but for sentences, where each sentence is processed by an LSTM

[Kiros et al., 2015]

One model to translate from any language to any other language





One model to translate from any language to any other language



Token designating target language



One model to translate from any language to any other language



Flipped order encoding

Token designating target language



One model to translate from any language to any other language



Flipped order encoding Why?

Token designating target language



One model to translate from any language to any other language



Flipped order encoding

Token designating target language

Johnson et al., 2017]

8! layer LSTM decoder and encoder



One model to translate from any language to any other language

8 layers Encoder LSTMs GPU8 **Residual** at other layers (ResNet style) GPU3 GPU2 GPU2 **Bi-directional** at lower layers GPU1 $\rightarrow \cdots \rightarrow$ Hello \rightarrow <2es> \rightarrow </s>

Flipped order encoding

Token designating target language



Johnson et al., 2017]

8! layer LSTM decoder and encoder



One model to translate from any language to any other language

8 layers Encoder LSTMs GPU8 **Residual** at other layers (ResNet style) GPU3 GPU2 GPU2 **Bi-directional** at lower layers GPU1 $\rightarrow \cdots \rightarrow$ Hello \rightarrow <2es> \rightarrow </s>

Flipped order encoding

Token designating target language



Johnson et al., 2017]

8! layer LSTM decoder and encoder





* slide from Dhruv Batra

, Dotr

Image Embedding (VGGNet)



* slide from Dhruv Batra

Image Embedding (VGGNet)



* slide from Dhruv Batra



Image Embedding (VGGNet)

* slide from Dhruv Batra

, Dotr



Image Embedding (VGGNet)

* slide from Dhruv Batra

, Dotr

Applications: Neural Image Captioning Good results



A cat sitting on a suitcase on the floor



Two people walking on the beach with surfboards



A cat is sitting on a tree branch



A tennis player in action on the court



A dog is running in the grass with a frisbee



Two giraffes standing in a grassy field



A white teddy bear sitting in the grass



A man riding a dirt bike on a dirt track

Applications: Neural Image Captioning Failure cases



A woman is holding a cat in her hand



A person holding a computer mouse on a desk



A woman standing on a beach holding a surfboard



A bird is perched on a tree branch



A man in a baseball uniform throwing a ball

RNN focuses its attention at a different spatial location when generating each word



[Xu et al., ICML 2015]





[Xu et al., ICML 2015]





[Xu et al., ICML 2015]





[Xu et al., ICML 2015]





[Xu et al., ICML 2015]





[Xu et al., ICML 2015]







[Xu et al., ICML 2015]





[Xu et al., ICML 2015]





[Xu et al., ICML 2015]



Applications: Image Captioning with Attention **Good** results



A woman is throwing a frisbee in a park.





A little girl sitting on a bed with a teddy bear.

A group of people sitting on a boat in the water.

[Xu et al., ICML 2015]

A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.





A giraffe standing in a forest with trees in the background.



Applications: Image Captioning with Attention Failure results



A large white bird standing in a forest.



A woman holding a clock in her hand.



A person is standing on a beach with a surfboard.

A woman is sitting at a table with a large pizza.

[Xu et al., ICML 2015]

A man wearing a hat and a hat on a skateboard.





A man is talking on his cell phone while another man watches.



Image



Question

"How many horses are in this image?"

* slide from Dhruv Batra

, Dotr

Image Embedding (VGGNet)



Question

"How many horses are in this

s image?"

* slide from Dhruv Batra

, Dotr

Image Embedding (VGGNet)



Question Embedding (LSTM)



* slide from Dhruv Batra

Image Embedding (VGGNet)





* slide from Dhruv Batra



Activity: A collection of human/object movements with a particular semantic meaning



Activity: A collection of human/object movements with a particular semantic meaning



Action Recognition: Finding if a video segment contains such a movement



Action Recognition: Finding if a video segment contains such a movement

Activity: A collection of human/object movements with a particular semantic meaning

Action Detection: Finding a segment (beginning and start) and recognize the action in it




Early Detection: Recognize when an action starts and try to predict which action is performed as quickly as possible.



video

Applications: Activity Detection Penalty at every time step is the same



video



Applications: Activity Detection Penalty at every time step is the same





- Detecting the correct action class
- More confident that it is not the incorrect action class



As the detector sees more of an action, it should become more confident of

- Detecting the correct action class
- More confident that it is not the incorrect action class



As the detector sees more of an action, it should become more confident of



- Detecting the correct action class
- More confident that it is not the incorrect action class



As the detector sees more of an action, it should become more confident of

- Detecting the correct action class
- More confident that it is not the incorrect action class



As the detector sees more of an action, it should become more confident of

cooking

New Class of Loss Functions

Training loss at time t: $\mathcal{L}^t =$

- \mathcal{L}_r^t is one of the following:

Classification loss at time t

$$\mathcal{L}_c^t + \lambda_r \mathcal{L}_r^t$$

Ranking loss at time t

• \mathcal{L}_s^t ranking loss on detection score • \mathcal{L}_m^t ranking loss on discriminative margin

Ideally what we want:



Prediction score of the ground truth action label



Prediction score of the ground truth action label



Prediction score of the ground truth action label





Prediction score of the ground truth action label

Activity detection performance measured in mAP at different IOU thresholds

Model	$\alpha = 0.1$	$\alpha = 0.2$	$\alpha = 0.3$	α = 0.4	$\alpha = 0.5$	α = 0.6	$\alpha = 0.7$	$\alpha = 0.8$
Heilbron <i>et al</i> .	12.5%	11.9%	11.1%	10.4%	9.7%	-	-	-
CNN	30.1%	26.9%	23.4%	21.2%	18.9%	17.5%	16.5%	15.8%
LSTM	48.1%	44.3%	40.6%	35.6%	31.3%	28.3%	26.0%	24.6%
LSTM-m	52.6%	48.9%	45.1%	40.1%	35.1%	31.8%	29.1%	27.2%
LSTM-s	54.0%	50.1%	46.3%	41.2%	36.4%	33.0%	30.4%	28.7%

LSTM-m LSTM trained using both classification loss and rank loss on *discriminative margin*. LSTM trained using both classification loss and rank loss on *detection score*. LSTM-s

Activity early detection performance measured in mAP at different IOU thresholds

Model	$\alpha = 0.1$	$\alpha = 0.2$	$\alpha = 0.3$	$\alpha = 0.4$	$\alpha = 0.5$	$\alpha = 0.6$	$\alpha = 0.7$	$\alpha = 0.8$
CNN	27.0%	23.4%	20.4%	17.2%	14.6%	12.3%	11.0%	10.3%
LSTM	49.5%	44.7%	38.8%	33.9%	29.6%	25.6%	23.5%	22.4%
LSTM-m	52.6%	47.9%	41.5%	36.2%	31.4%	27.1%	24.8%	23.5%
LSTM-s	55.1%	50.3%	44.0%	38.9%	34.1%	29.8%	27.4%	26.1%

LSTM-m LSTM trained using both classification loss and rank loss on *discriminative margin*. LSTM trained using both classification loss and rank loss on *detection score*. LSTM-s

Note: first 3/10 of activity is seen by a detector

Activity early detection performance measured in mAP at different IOU thresholds

Model	$\alpha = 0.1$	$\alpha = 0.2$	$\alpha = 0.3$	$\alpha = 0.4$	$\alpha = 0.5$	$\alpha = 0.6$	$\alpha = 0.7$	$\alpha = 0.8$
CNN	27.0%	23.4%	20.4%	17.2%	14.6%	12.3%	11.0%	10.3%
LSTM	49.5%	44.7%	38.8%	33.9%	29.6%	25.6%	23.5%	22.4%
LSTM-m	52.6%	47.9%	41.5%	36.2%	31.4%	27.1%	24.8%	23.5%
LSTM-s	55.1%	50.3%	44.0%	38.9%	34.1%	29.8%	27.4%	26.1%

LSTM-m LSTM trained using both classification loss and rank loss on *discriminative margin*. LSTM trained using both classification loss and rank loss on *detection score*. LSTM-s

Take home: Early detection is only 1-3% worse than sewing the whole sequence

Note: first 3/10 of activity is seen by a detector









Attention Models for Action Highlighting





(b) Action: Having Massage

[Torabi & Sigal, 2017]

2017