

THE UNIVERSITY OF BRITISH COLUMBIA

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

Lecture 10: RNN Applications



Course Logistics

- Assignment 3 due date is Wednsday 11:59pm (Thursday?)
- Assignment 4 is out, due Friday, February 15th @ 11:59pm

- Group **Projects** form completed - Project proposals (in class on **Feb 26th**)

Final **Project** (50% of grade total) – Reminder

- Group project (groups of 3 are encouraged, but fewer maybe possible)
- Groups are self-formed, you will not be assigned to a group
- You need to come up with a project proposal and then work on the project as a group (each person in the group gets the same grade for the project)
- Project needs to be research oriented (not simply implementing an existing) paper); you can use code of existing paper as a starting point though

Project proposal + class presentation: 15% Project + final presentation: 35%

Project proposal and class presentation – 15% of grade

Presentation (~5 minutes irrespective of the group size)

- 1. Clear explanation of the overall problem you want to solve and relationship to the topics covered in class
- 2. What **model/algorithms** you planning to explore: at this can be somewhat abstract (e.g., CNN+RNN)
- 3. The **dataset(s)** you will use and how will you **evaluate** performance
- 4. List of **papers** you plan to read as references
- 5. How will you structure the project, who will do what and a rough timeline

Project proposal and class presentation – 15% of grade

Presentation (~5 minutes irrespective of the group size)

- 2. What **model/algorithms** you planning to explore: this can be somewhat abstract (e.g., CNN+RNN)
- 3. The **dataset(s)** you will use and how will you **evaluate** performance
- 4. List of **papers** you plan to read as references
- 5. How will you structure the project, who will do what and a rough timeline

After presentation you will get the feedback from me

1. Clear explanation of the overall problem you want to solve and relationship to the topics covered in class

Project proposal and class presentation – 15% of grade

Presentation (~5 minutes irrespective of the group size)

- 1. Clear explanation of the overall problem you want to solve and relationship to the topics covered in class
- 2. What **model/algorithms** you planning to explore: this can be somewhat abstract (e.g., CNN+RNN)
- 3. The **dataset(s)** you will use and how will you **evaluate** performance
- 4. List of **papers** you plan to read as references
- 5. How will you structure the project, who will do what and a rough timeline

After presentation you will get the feedback from me

Proposal

- Same as above but in more detail, with well defined algorithms and timeline
- Will be in the form of the **PDF** document (initial paper draft)

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)





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]

* slide from Marco Pedersoli and Thomas Lucas

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



[Johnson et al., 2017]



One model to translate from any language to any other language



Token designating target language

[Johnson et al., 2017]



One model to translate from any language to any other language



Flipped order encoding Why?

Token designating target language

[Johnson et al., 2017]



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

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



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

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]