



Language Technologies Institute



## **Multimodal Machine Learning**

## Lecture 3.2: Recurrent Networks Louis-Philippe Morency

\* Original version co-developed with Tadas Baltrusaitis

#### **Lecture Objectives**

- Sequential modeling with convolutional networks
- Word representations
  - Distributional hypothesis
  - Learning neural representations
- Language models and sequence modeling tasks
- Recurrent neural networks
  - Gated recurrent neural networks
  - Long Short-Term Memory (LSTM) model
    - Multi-view LSTM
  - Backpropagation through time



# Sequential Modeling with Convolutional Networks



Language Technologies Institute

## **Modeling Temporal and Sequential Data**



#### How to represent a video sequence?

#### One option: Recurrent Neural Networks (more about this on Thursday)



## **3D CNN**



## Input as a 3D tensor (stacking video images)

**3D CNN** 



First layer with 3D kernels



## **Time-Delay Neural Network**



**Alexander Waibel**, Phoneme Recognition Using Time-Delay Neural Networks, SP87-100, Meeting of the Institute of Electrical, Information and Communication Engineers (IEICE), December, 1987, Tokyo, Japan.



#### Temporal Convolution Network (TCN) [Lea et al., CVPR 2017]

#### RRRRRRRRRRR





### Dilated TCN Model [Lea et al., CVPR 2017]

#### **Dilated Convolutions**



Dilation of 4: Step size of 4 when convoluting

#### + Skip connections to help with deep modeling



#### Dilated TCN Models [Lea et al., CVPR 2017]







# Representing Words: Distributed Semantics



Language Technologies Institute

## What is the meaning of "bardiwac"?

- He handed her glass of bardiwac.
- Beef dishes are made to complement the bardiwacs.
- Nigel staggered to his feet, face flushed from too much bardiwac.
- Malbec, one of the lesser-known bardiwac grapes, responds well to Australia's sunshine.
- I dined off bread and cheese and this excellent bardiwac.
- The drinks were delicious: blood-red bardiwac as well as light, sweet Rhenish.
- ⇒ bardiwac is a heavy red alcoholic beverage made from grapes



## **The Distributional Hypothesis**

- Distribution Hypothesis (DH) [Lenci 2008]
  - At least certain aspects of the meaning of lexical expressions depend on their distributional properties in the linguistic contexts
  - The degree of semantic similarity between two linguistic expressions  $\alpha$  and  $\beta$  is a function of the similarity of the linguistic contexts in which  $\alpha$  and  $\beta$  can appear
- Weak and strong DH
  - Weak view as a quantitative method for semantic analysis and lexical resource induction
  - Strong view as a cognitive hypothesis about the form and origin of semantic representations; assuming that word distributions in context play a specific *causal role* in forming meaning representations.



### **Geometric interpretation**

- row vector X<sub>dog</sub> describes usage of word *dog* in the corpus
- can be seen as coordinates of point in *n*-dimensional Euclidean space R<sup>n</sup>

	get	see	use	hear	eat	kill
knife	51	20	84	0	3	0
cat	52	58	4	4	6	26
dog	115	83	10	42	33	17
boat	59	39	23	4	0	0
cup	98	14	6	2	1	0
pig	12	17	3	2	9	27
banana	11	2	2	0	18	0

#### co-occurrence matrix M



## **Distance and similarity**

- illustrated for two dimensions: get and use: X<sub>dog</sub> = (115, 10)
- similarity = spatial proximity (Euclidean distance)

nse

■ location depends on frequency of noun  $(f_{dog} \approx 2.7 \cdot f_{cat})$  Two dimensions of English V-Obj DSM





## Angle and similarity

- direction more important than location
- normalise "length"
   ||x<sub>dog</sub>|| of vector
- or use angle α as distance measure

#### Two dimensions of English V-Obj DSM



use

get



## **Semantic maps**





# Learning Neural Word Representations



Language Technologies Institute

#### How to learn neural word representations?

- Distribution hypothesis: Approximate the word meaning by its surrounding words
- Words used in a similar context will lie close together





Instead of capturing co-occurrence counts directly, predict surrounding words of every word

$$\frac{1}{T} \sum_{t=1}^{T} \sum_{-c \le j \le c, j \ne 0} \log p(w_{t+j} | w_t)$$



#### How to learn neural word representations?



#### Language Technologies Institute

#### How to use these word representations

If we would have a vocabulary of 100 000 words:





#### **Vector space models of words**

- While learning these word representations, we are actually building a vector space in which all words reside with certain relationships between them
- Encodes both syntactic and semantic relationships



This vector space allows for algebraic operations:

Vec(king) – vec(man) + vec(woman) ≈ vec(queen)

Why linear algebra is working?



#### **Vector space models of words: semantic relationships**



Trained on the Google news corpus with over 300 billion words



# Language Sequence Modeling Tasks



Language Technologies Institute

## **Sequence Modeling: Sequence Label Prediction**



By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful



Sentiment ? (positive or negative)





## **Sequence Modeling: Sequence Prediction**



By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful



Part-of-speech ? (noun, verb,...)





## **Sequence Modeling: Sequence Representation**







## **Sequence Modeling: Language Model**



By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful



#### Language Model





## **Application: Speech Recognition**

arg max P(wordsequence | acoustics) = wordsequence

$$\underset{wordsequence}{\operatorname{arg\,max}} \frac{P(acoustics \mid wordsequence) \times P(wordsequence)}{P(acoustics)}$$

 $arg \max P(acoustics | wordsequence) \times P(wordsequence)$ 

wordsequence







## **Application: Language Generation**



Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

#### Example: Image captioning





## **N-Gram Language Model Formulations**

- Word sequences  $w_1^n = w_1 \dots w_n$
- Chain rule of probability  $P(w_1^n) = P(w_1)P(w_2 \mid w_1)P(w_3 \mid w_1^2)...P(w_n \mid w_1^{n-1}) = \prod_{k=1}^n P(w_k \mid w_1^{k-1})$
- Bigram approximation
   P(w<sub>1</sub><sup>n</sup>) = \product P(w<sub>k</sub> | w<sub>k-1</sub>)
   N-gram approximation
  - $P(w_1^n) = \prod_{k=1}^n P(w_k \mid w_{k-N+1}^{k-1})$



## **Evaluating Language Model: Perplexity**

The best language model is one that best predicts an unseen test set

Chain rule:

For bigrams:

• Gives the highest P(sentence)

Perplexity is the inverse probability of the test set, normalized by the number of words:

$$PP(W) = P(w_1w_2...w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1w_2...w_N)}}$$

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_1\dots w_{i-1})}}$$

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_{i-1})}}$$



## **Challenges in Sequence Modeling**



- Part-of-speech ? (noun, verb,...)
- Sentiment ? (positive or negative)
- Language Model
- Sequence representation

## Main Challenges:

- Sequences of variable lengths (e.g., sentences)
- Keep the number of parameters at a minimum
- Take advantage of possible redundancy



Language Technologies Institute

# Recurrent Neural Networks



Language Technologies Institute

## **Recurrent Neural Network**

#### **Feedforward Neural Network**







#### **Recurrent Neural Networks**





#### **Recurrent Neural Networks - Unrolling**



Same model parameters are used for all time parts.





#### **RNN-based Language Model**



Models long-term information

#### **RNN-based Sentence Generation (Decoder)**



Models long-term information

## **Sequence Modeling: Sequence Prediction**



By Antony Witheyman - January 12, 2006

Ideal for anyone with an interest in disguises who likes to see the subject tackled in a humourous manner.

0 of 4 people found this review helpful



Sentiment ? (positive or negative)





## **RNN for Sequence Prediction**



$$L = \frac{1}{N} \sum_{t} L^{(t)} = \frac{1}{N} \sum_{t} -logP(Y = y^{(t)} | \mathbf{z}^{(t)})$$

### **RNN for Sequence Prediction**



 $L = L^{(N)} = -logP(Y = y^{(N)} | \mathbf{z}^{(N)})$ 

## **Sequence Modeling: Sequence Representation**







#### **RNN for Sequence Representation**



## **RNN for Sequence Representation (Encoder)**



**RNN-based for Machine Translation** 

Le chien sur la plage



The dog on the beach



#### **Encoder-Decoder Architecture**



## **Related Topics**

- Character-level "language models"
  - Xiang Zhang, Junbo Zhao and Yann LeCun, Character-level Convolutional Networks for Text Classification, NIPS 2015

http://arxiv.org/pdf/1509.01626v2.pdf

Skip-though: embedding at the sentence level

 Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Antonio Torralba, Raquel Urtasun, Sanja Fidler. Skip-Thought Vectors, NIPS 2015

http://arxiv.org/pdf/1506.06726v1.pdf



# Gated Recurrent Neural Networks



Language Technologies Institute



## **Long-term Dependencies**

Vanishing gradient problem for RNNs:



The influence of a given input on the hidden layer, and therefore on the network output, either decays or blows up exponentially as it cycles around the network's recurrent connections.

49



### **Recurrent Neural Networks**







## LSTM ideas: (1) "Memory" Cell and Self Loop

[Hochreiter and Schmidhuber, 1997]

Long Short-Term Memory (LSTM)





51 ogies Institute



## LSTM Ideas: (2) Input and Output Gates

[Hochreiter and Schmidhuber, 1997]





Language Technologies Institute

## LSTM Ideas: (3) Forget Gate [Gers et al., 2000]





#### **Recurrent Neural Network using LSTM Units**



Gradient can still be computer using backpropagation!



#### **Bi-directional LSTM Network**







### **Deep LSTM Network**





**Carnegie Mellon University** 

56

#### **Multimodal Sequence Modeling – Early Fusion**







## Multi-View Long Short-Term Memory (MV-LSTM)







## **Multi-View Long Short-Term Memory**





## **Topologies for Multi-View LSTM**







## Multi-View Long Short-Term Memory (MV-LSTM)

#### Multimodal prediction of children engagement

Class labels	Model	Precision	Recall	F1
Easy to engage	LSTM (Early fusion)	0.75	0.81	0.78
	MV-LSTM Full	0.81	0.81	0.81
	MV-LSTM Coupled	0.79	0.81	0.80
	MV-LSTM Hybrid	0.80	0.86	0.83
Difficult to engage	LSTM (Early fusion)	0.63	0.55	0.59
	MV-LSTM Full	0.68	0.68	0.68
	MV-LSTM Coupled	0.67	0.64	0.65
	MV-LSTM Hybrid	0.74	0.64	0.68





# Backpropagation Through Time



Language Technologies Institute

## **Optimization: Gradient Computation**

#### Vector representation:





Language Technologies Institute

## **Backpropagation Algorithm**

### Forward pass

 Following the graph topology, compute value of each unit

### **Backpropagation pass**

- Initialize output gradient = 1
- Compute "local" Jacobian matrix using values from forward pass
- Use the chain rule:

```
Gradient = "local" Jacobian x
"backprop" gradient
```





#### **Recurrent Neural Networks**





## **Backpropagation Through Time**

$$L = \sum_{t} L^{(t)} = -\sum_{t} log P(Y = y^{(t)} | z^{(t)})$$

$$(L^{(t)} \text{ or } L^{(t)}) \frac{\partial L}{\partial L^{(t)}} = 1$$

$$(T^{(t)} \text{ or } L^{(t)}) \frac{\partial L}{\partial L^{(t)}} = 1$$

$$(T^{(t)} \text{ or } Z^{(t)}) \frac{\partial L}{\partial L^{(t)}} = \frac{\partial L}{\partial z_{i}^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial z_{i}^{(t)}} = sigmoid(z_{i}^{t}) - \mathbf{1}_{i,y^{(t)}}$$

$$(T^{(t)} P_{h^{(t)}}L = P_{z^{(t)}}L \frac{\partial z^{(t)}}{\partial h^{(t)}} = P_{z^{(t)}}LV$$

$$(T^{(t)} P_{h^{(t)}}L = P_{z^{(t)}}L \frac{\partial o^{(t)}}{\partial h^{(t)}} + P_{z^{(t+1)}}L \frac{\partial h^{(t+1)}}{\partial h^{(t)}}$$



Language Technologies Institute

#### **Carnegie Mellon University**

 $\tau$ )

 $\mathbf{z}^{(\tau)}$ 

 $h^{( au)}$ 

 $\mathbf{x}^{(\tau)}$ 

## **Backpropagation Through Time**

$$L = \sum_{t} L^{(t)} = -\sum_{t} log P(Y = y^{(t)} | \mathbf{z}^{(t)})$$
  
Gradient = "backprop" gradient  
x "local" Jacobian

$$\bigcup \nabla_{\boldsymbol{U}} L = \sum_{t} (\nabla_{\boldsymbol{h}^{(t)}} L) \frac{\partial \boldsymbol{h}^{(t)}}{\partial \boldsymbol{U}}$$



