

#### 496/896 Lecture 6: Recurrent Architectures Stephen Scott Introduction Basic Idea I/O Mappings Examples Training Deep RNNs

GRUs

# Additional CSCE 496/896 Lecture 6: Recurrent Architectures Stephen Scott (Adapted from Vinod Variyam and Ian Goodfellow)

## sscott@cse.unl.edu

# Nebraska Introduction



GRUs

- All our architectures so far work on fixed-sized inputs
- Recurrent neural networks work on sequences of inputs
- E.g., text, biological sequences, video, audio
- Can also try 1D convolutions, but lose long-term relationships in input
- Especially useful for NLP applications: translation, speech-to-text, sentiment analysis
- Can also **create novel output:** e.g., Shakespearean text, music

・ロト・西ト・モー・ 一日・ うへの





# Nebraska Basic Recurrent Layer

#### 6 6: nt res

Basic Idea

Examples

Deep RNNs

GRUs

I/O Mappings

- Each node in the recurrent layer has independent weights for both  $\mathbf{x}_{(t)}$  and  $\mathbf{y}_{(t-1)}$
- For a single recurrent node, denote by  $w_x$  and  $w_y$
- For the entire layer, combine into matrices  $\mathit{W}_x$  and  $\mathit{W}_y$
- For activation function  $\phi$  and bias vector  $\pmb{b},$  output vector is

$$\mathbf{y}_{(t)} = \phi \left( W_{\mathbf{x}}^{\top} \mathbf{x}_{(t)} + W_{\mathbf{y}}^{\top} \mathbf{y}_{(t-1)} + \mathbf{b} \right)$$

#### Nebraska Memory and State

- Since a node's output depends on its past, it can be thought of having memory or state
- State at time *t* is  $\boldsymbol{h}_{(t)} = f(\boldsymbol{h}_{(t-1)}, \boldsymbol{x}_{(t)})$  and output  $\boldsymbol{y}_{(t)} = g(\boldsymbol{h}_{(t-1)}, \boldsymbol{x}_{(t)})$
- State could be the same as the output, or separate • Can think of  $h_{(t)}$  as storing important information about
- input sequence Analogous to convolutional outputs summarizing important image features



#### Input/Output Mappings Nebraska Sequence to Sequence

96/896 ecture 6

tenhen Sr

Basic Idea I/O Mappings

Examples

Training

GRUs

Deep RNNs

ntroduction

I/O Mappings

Examples

Deep RNNs

Basic Idea

I/O Mappings

Examples

Deep RNNs

GRUs

э

Training

LSTMs

Many ways to employ this basic architecture:

- Sequence to sequence: Input is a sequence and output is a sequence
- E.g., series of stock predictions, one day in advance



#### Input/Output Mappings Nebraska Vector to Sequence

- Vector to sequence: Input is a single vector (zeroes) for other times) and output is a sequence
  - E.g., image to caption



#### Input/Output Mappings Nebraska Encoder-Decoder Architecture

(0)

- Encoder-decoder: Sequence-to-vector (encoder) followed by vector-to-sequence (decoder)
- Input sequence  $(x_1, \ldots, x_T)$  yields hidden outputs  $(h_1, \ldots, h_T)$ , then mapped to **context vector**  $\boldsymbol{c} = f(\boldsymbol{h}_1, \ldots, \boldsymbol{h}_T)$
- Decoder output y<sub>t</sub> depends on previously output  $(y_1, ..., y_{t'-1})$  and *c*
- Example application: neural machine translation



#### Input/Output Mappings Nebraska Encoder-Decoder Architecture: NMT Example

- Pre-trained word embeddings fed into input
- Encoder maps word sequence to vector, decoder maps to translation via softmax distribution
- After training, do translation by feeding previous translated word  $y'_{(t-1)}$  to decoder



# Basic Idea O Mapping Examples Training Deep RNNs GRUs

Nebraska

ntroductior

xamples

Deep RNNs

Stephen Sco

Basic Idea

I/O Mappings

Examples

Training Deep RNNs

GBUs

Training

LSTMs

## Input/Output Mappings Sequence to Vector

• Sequence to vector: Input is sequence and output a vector/score/ classification



Х

(1)

(2)

(3) **X**<sub>(3)</sub> х

# I/O Mappings

### Input/Output Mappings Encoder-Decoder Architecture

- Works through an embedded space like an autoencoder, so can represent the entire input as an embedded vector prior to decoding
- Issue: Need to ensure that the context vector fed into decoder is sufficiently large in dimension to represent context required
- Can address this representation problem via attention mechanism mechanism
  - Encodes input sequence into a vector sequence rather than single vector
  - As it decodes translation, decoder focuses on relevant subset of the vectors

#### Netranska Linola E-D Architecture: Attention Mechanism (Bahdanau et al., 2015)

 Bidirectional RNN reads input forward and backward simultaneously

Stenhen So

troduction

Basic Idea

I/O Mappings

Examples

Deep RNNs

GRUs

 Encoder builds annotation h<sub>j</sub> as concatenation of h
<sub>j</sub> and h
<sub>j</sub>
 ⇒ h<sub>j</sub> summarizes preceding and following inputs
 ith context vector

$$c_{i} = \sum_{j=1}^{T} \alpha_{ij} h_{j}, \text{ where}$$
$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{i=1}^{T} \exp(e_{ik})}$$



・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

and *e<sub>ij</sub>* is an **alignment score** between inputs around *j* and outputs around *i* 

# Nebraska

ntroduction

I/O Mappings

Examples Training

Deep RNNs

LSTMs

Nebraska

Basic Idea

Examples

Deep RNNs

Training

GRUs

I/O Mappings

Input/Output Mappings E-D Architecture: Attention Mechanism (Bahdanau et al., 2015)

• The *i*th element of **attention vector**  $\alpha_j$  tells us the probability that target output  $y_i$  is aligned to (or translated from) input  $x_j$ 



- Then c<sub>i</sub> is expected annotation over all annotations with probabilities α<sub>i</sub>
- Alignment score e<sub>ij</sub> indicates how much we should focus on word encoding h<sub>j</sub> when generating output y<sub>i</sub> (in decoder state s<sub>i-1</sub>)
- Can compute  $e_{ij}$  via dot product  $h_j^{\top} s_{i-1}$ , bilinear function  $h_i^{\top} W s_{i-1}$ , or nonlinear activation  $a_{i-1} + a_{i-1} + a_{$

# Nebraska Linoh Static Unrolling for Two Time Steps

CSCE 496/896 Lecture 6: Recurrent Architectures	
Stephen Scott	<pre>X0 = tf.placeholder(tf.float32, [None, n_inputs]) X1 = tf.placeholder(tf.float32, [None, n_inputs]) Wx = tf.Variable(tf.random_normal(shape=[n_inputs, n_neurons],dtype=tf.float32))</pre>
ntroduction	<pre>Wy = tf.Variable(tf.random_normal(shape=[n_neurons,n_neurons],dtype=tf.float32)) b = tf.Variable(tf.zeros([1, n_neurons], dtype=tf.float32))</pre>
asic Idea O Mappings	Y0 = tf.tanh(tf.matmul(X0, Wx) + b) Y1 = tf.tanh(tf.matmul(Y0, Wy) + tf.matmul(X1, Wx) + b)
xamples	Input:
Deep RNNs STMs GRUs	<pre># Mini-batch: instance 0, instance 1, instance 2, instance 3 X0_batch = np.array([[0, 1, 2], [3, 4, 5], [6, 7, 8], [9, 0, 1]]) # t = 0 X1_batch = np.array([[9, 8, 7], [0, 0, 0], [6, 5, 4], [3, 2, 1]]) # t = 1</pre>

Nebiaska	Example Implementation	Nebraska	Example Implementation
Lincoln	Static Unrolling for Two Time Steps	Lincoln	Automatic Static Unrolling
CSCE 496/896 Lecture 6: Recurrent Architectures Stephen Scott Introduction Basic Idea I/O Mappings Examples Training Deep RINIs LSTMs GRUS	Can achieve the same thing more compactly via static_rnn() x0 = tf.placeholder(tf.float32, [None, n_inputs]) x1 = tf.placeholder(tf.float32, [None, n_inputs]) basic_cell = tf.contrib.rnn.statio_rnn(basic_cell, [X0, x1], dtype=tf.float32) y0, y1 = output_seqs Automatically unrolls into length-2 sequence RNN	CSCE 496/896 Lecture 6: Recurrent Architectures Stephen Scott Introduction Basic Idea I/O Mappings Examples Training Deep RNNs LSTMs GRUs	<pre>Can avoid specifying one placeholder per time step via tf.stack and tf.unstack  X = tf.placeholder(tf.float32, [None, n.steps, n_inputs])     X_seqs = tf.unstack(tf.transpose(X, perm=[1, 0, 2]))     basic_cell = tf.contrib.rn.BasicRNNCell(nm_units=n_neurons)     outputs=gs, states = tf.contrib.rn.state_rn(basic_cell, X_seqs,</pre>



Training Nebraska Nebraska Backpropagation Through Time (BPTT) CSCE Unroll through time and use BPTT 496/896 Lecture 6 Recurrent ecture 6 • Forward pass mini-batch of sequences through unrolled Recurren Inchitectur Architectur network yields output sequence  $Y_{(t_{\min})}, \ldots, Y_{(t_{\max})}$  Output sequence evaluated using cost ntroduction ntroductior  $C\left(Y_{(t_{\min})},\ldots,Y_{(t_{\max})}\right)$ asic Idea asic Idea Gradients propagated backward through unrolled I/O Mappings I/O Mapping network (summing over all time steps), and parameters Examples Examples  $C(\mathbf{Y}_{(2)}, \mathbf{Y}_{(3)}, \mathbf{Y}_{(4)})$ Training Deep RNNs Deep RNNs I STMs **Y**<sub>(2)</sub> Y<sub>(4)</sub> Υ<sub>(3)</sub> À Á W,b W,b W,b W,b W,b **X**<sub>(4)</sub> ∃ • • • • • 3.5

Training

I STMs

・ロト・西ト・ヨト・ヨト・ヨー もくの

Training Example: Training on MNIST as a Vector Sequence

- Consider MNIST inputs provided as sequence of 28 inputs of 28-dimensional vectors
  - Feed in input as usual, then compute loss between target and softmax output after 28th input





# Nebraska

Stenhen Sco

ntroduction Basic Idea

O Mapping Examples

Training

Deep RNNs GRUs

Training Example: Training on Time Series Data

- Use sequences of length n\_steps=20 and n\_neurons=100 recurrent neurons
- Since output size = 100 > 1 = target size, use OutputProjectionWrapper to feed recurrent layer output into a linear unit to get a scalar



Nebraska	Training
Lincoln	Example: Training on Time Series Data
CSCE 496/896 Lecture 6: Recurrent Architectures Stephen Scott Introduction Basic Idea I/O Mappings Examples Training Deep RNNs LSTMs GRUs 26/35	<pre>n_steps = 20 n_inputs = 1 n_outputs = 1 2 = f_placeholder(ff.float32, [None, n_steps, n_inputs]) y = f.f_placeholder(ff.float32, [None, n_steps, n_outputs]) cell = (f_oontrib.rnn.duptprejectionKrapp(</pre>





Lincoin	
CSCE 496/896 e.ecture 6: Recurrent chitectures ephen Scott oduction	<ul> <li>Vanishing and exploding gradients can be a problem with RNNs, like with other deep networks</li> <li>Can as usual address with, e.g., ReLU, batch normalization, gradient clipping, etc.</li> <li>Can still suffer from long training times with long input sequences</li> </ul>
Mappings amples ining	<ul> <li>Truncated backpropagation through time addresses this by limiting n_steps</li> <li>Lose ability to learn long-term patterns</li> </ul>
ep RNNs T <mark>Ms</mark> Us	<ul> <li>In general, also have problem of first inputs of sequence have diminishing impact as sequence grows</li> <li>E.g., first few words of long text sequence</li> <li>Goal: Introduce long-term memory to RNNs</li> </ul>
	• Allow a network to accumulate information about the

(E) (E) = 9000 < @ >





GRU ce

