NPFL114, Lecture 7



Recurrent Neural Networks

Milan Straka

🖬 April 14, 2020





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated



Recurrent Neural Networks

NPFL114, Lecture 7

GRU

HighwayNetworks

RNNRegularization

RNNApplications

Recurrent Neural Networks



Single RNN cell



Unrolled RNN cells







Given an input $m{x}^{(t)}$ and previous state $m{s}^{(t-1)}$, the new state is computed as

$$oldsymbol{s}^{(t)} = f(oldsymbol{s}^{(t-1)},oldsymbol{x}^{(t)};oldsymbol{ heta}).$$

One of the simplest possibilities (called SimpleRNN in TensorFlow) is

$$oldsymbol{s}^{(t)} = anh(oldsymbol{U}oldsymbol{s}^{(t-1)} + oldsymbol{V}oldsymbol{x}^{(t)} + oldsymbol{b}).$$

Basic RNN Cell

Basic RNN cells suffer a lot from vanishing/exploding gradients (*the challenge of long-term dependencies*).

If we simplify the recurrence of states to

$$oldsymbol{s}^{(t)} = oldsymbol{U}oldsymbol{s}^{(t-1)},$$

we get

$$oldsymbol{s}^{(t)} = oldsymbol{U}^t oldsymbol{s}^{(0)}.$$

If U has eigenvalue decomposition of $oldsymbol{U} = oldsymbol{Q}oldsymbol{\Lambda}oldsymbol{Q}^{-1}$, we get

$$oldsymbol{s}^{(t)} = oldsymbol{Q}oldsymbol{\Lambda}^toldsymbol{Q}^{-1}oldsymbol{s}^{(0)}.$$

The main problem is that the *same* function is iteratively applied many times.

Several more complex RNN cell variants have been proposed, which alleviate this issue to some degree, namely **LSTM** and **GRU**.

NPFL114, Lecture 7

LSTM GRU

RNN

HighwayNetworks

RNNRegularization

RNNApplications



Hochreiter & Schmidhuber (1997) suggested that to enforce *constant error flow*, we would like

$$f'=\mathbf{1}.$$

They propose to achieve that by a *constant error carrousel*.



They also propose an *input* and *output* gates which control the flow of information into and out of the carrousel (*memory cell* c_t).

$$egin{aligned} oldsymbol{i}_t &\leftarrow \sigma(oldsymbol{W}^ioldsymbol{x}_t + oldsymbol{V}^ioldsymbol{h}_{t-1} + oldsymbol{b}^i) \ oldsymbol{o}_t &\leftarrow \sigma(oldsymbol{W}^ooldsymbol{x}_t + oldsymbol{V}^ooldsymbol{h}_{t-1} + oldsymbol{b}^o) \ oldsymbol{c}_t &\leftarrow oldsymbol{c}_{t-1} + oldsymbol{i}_t \cdot anh(oldsymbol{W}^yoldsymbol{x}_t + oldsymbol{V}^yoldsymbol{h}_{t-1} + oldsymbol{b}^y) \ oldsymbol{h}_t &\leftarrow oldsymbol{o}_t \cdot anh(oldsymbol{c}_t) \end{aligned}$$



GRU

RNNRegularization



Later in Gers, Schmidhuber & Cummins (1999) a possibility to *forget* information from memory cell c_t was added.

$$egin{aligned} egin{aligned} egi$$

GRU

RNNRegularization

RNNApplications

WordEmbeddings







GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings





NPFL114, Lecture 7

LSTM RNN

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings





http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-C-line.png

NPFL114, Lecture 7

RNN LSTM

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings





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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-i.png

NPFL114, Lecture 7

RNN LSTM

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings





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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-f.png

NPFL114, Lecture 7

RNN LSTM

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings





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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-C.png

NPFL114, Lecture 7

RNN LSTM

HighwayNetworks

GRU

RNNRegularization

RNNApplications

WordEmbeddings





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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-focus-o.png

NPFL114, Lecture 7

RNN LSTM

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings

Gated Recurrent Unit

Gated recurrent unit (GRU) was proposed by Cho et al. (2014) as a simplification of LSTM. The main differences are

- no memory cell
- forgetting and updating tied together

$$egin{aligned} m{r}_t &\leftarrow \sigma(m{W}^rm{x}_t + m{V}^rm{h}_{t-1} + m{b}^r) \ m{u}_t &\leftarrow \sigma(m{W}^um{x}_t + m{V}^um{h}_{t-1} + m{b}^u) \ m{\hat{h}}_t &\leftarrow anh(m{W}^hm{x}_t + m{V}^h(m{r}_t \cdot m{h}_{t-1}) + m{b}^h) \ m{h}_t &\leftarrow m{u}_t \cdot m{h}_{t-1} + (1 - m{u}_t) \cdot m{\hat{h}}_t \end{aligned}$$



GRU



Gated Recurrent Unit





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

http://colah.github.io/posts/2015-08-Understanding-LSTMs/img/LSTM3-var-GRU.png

NPFL114, Lecture 7

GRU

RNNRegularization

RNNApplications

WordEmbeddings



NPFL114, Lecture 7

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings



For input \boldsymbol{x} , fully connected layer computes

$$oldsymbol{y} \leftarrow H(oldsymbol{x},oldsymbol{W}_H).$$

Highway networks add residual connection with gating:

$$oldsymbol{y} \leftarrow H(oldsymbol{x},oldsymbol{W}_H) \cdot T(oldsymbol{x},oldsymbol{W}_T) + oldsymbol{x} \cdot (1 - T(oldsymbol{x},oldsymbol{W}_T)).$$

Usually, the gating is defined as

$$T(\boldsymbol{x}, \boldsymbol{W}_T) \leftarrow \sigma(\boldsymbol{W}_T \boldsymbol{x} + \boldsymbol{b}_T).$$

Note that the resulting update is very similar to a GRU cell with h_t removed; for a fully connected layer $H(x, W_H) = \tanh(W_H x + b_H)$ it is exactly it.

NPFL114, Lecture 7

GRU

RNNRegularization



Figure 1: Comparison of optimization of plain networks and highway networks of various depths. *Left:* The training curves for the best hyperparameter settings obtained for each network depth. *Right:* Mean performance of top 10 (out of 100) hyperparameter settings. Plain networks become much harder to optimize with increasing depth, while highway networks with up to 100 layers can still be optimized well. Best viewed on screen (larger version included in Supplementary Material).

NPFL114, Lecture 7RNNLSTMGRUHighwayNetworksRNNRegularizationRNNApplicationsWordEmbeddings20/35



NPFL114, Lecture 7

GRU

zation RNNApplications

Ú F_AL



Figure 4: Lesioned training set performance (y-axis) of the best 50-layer highway networks on MNIST (left) and CIFAR-100 (right), as a function of the lesioned layer (x-axis). Evaluated on the full training set while forcefully closing all the transform gates of a single layer at a time. The non-lesioned performance is indicated as a dashed line at the bottom.

Figure 4 of paper "Training Very Deep Networks", https://arxiv.org/abs/1507.06228.

Regularizing RNNs

Dropout

- Using dropout on hidden states interferes with long-term dependencies.
- However, using dropout on the inputs and outputs works well and is used frequently.
 In case residual connections are present, the output dropout needs to be applied before adding the residual connection.
- Several techniques were designed to allow using dropout on hidden states.
 - Variational Dropout
 - Recurrent Dropout
 - \circ Zoneout

GRU



Regularizing RNNs



Variational Dropout



To implement variational dropout on inputs in TensorFlow, use noise_shape of tf.keras.layers.Dropout to force the same mask across time-steps. The variational dropout on the hidden states can be implemented using recurrent_dropout argument of tf.keras.layers.{LSTM,GRU,SimpleRNN}{,Cell}.

NPFL114, Lecture 7

LSTM GRU

RNN

HighwayNetworks

RNNRegularization

RNNApplications

Recurrent Dropout

Dropout only candidate states (i.e., values added to the memory cell in LSTM and previous state in GRU).

Zoneout

Randomly preserve hidden activations instead of dropping them.

Batch Normalization

Very fragile and sensitive to proper initialization – there were papers with negative results (Dario Amodei et al, 2015: Deep Speech 2 or Cesar Laurent et al, 2016: Batch Normalized Recurrent Neural Networks) until people managed to make it work (Tim Cooijmans et al, 2016: Recurrent Batch Normalization; specifically, initializing $\gamma = 0.1$ did the trick).

GRU



variance causes vanishing gradient.



Figure 1 of paper "Recurrent Batch Normalization", https://arxiv.org/abs/1603.09025.

RNNRegularization

Layer Normalization

Much more stable than batch normalization.



Figure 2: Validation curves for the attentive reader model. BN results are taken from [Cooijmans et al., 2016].

Figure 2 of paper "Layer Normalization", https://arxiv.org/abs/1607.06450.

GRU

tion RNNApplications

ons WordEmbeddings

Layer Normalization

In an important recent architecture (namely Transformer), many fully connected layers are used, with a residual connection and a layer normalization.



This could be considered an alternative to highway networks, i.e., a suitable residual connection for fully connected layers. Note the architecture can be considered as a variant of a mobile inverted bottleneck 1×1 convolution block.

Basic RNN Applications



Sequence Element Representation

Create output for individual elements, for example for classification of the individual elements.



Sequence Representation

Generate a single output for the whole sequence (either the last output of the last state).

GRU

RNNRegularization



Sequence Prediction

During training, predict next sequence element.



During inference, use predicted elements as further inputs.



Multilayer RNNs

We might stack several layers of recurrent neural networks. Usually using two or three layers gives better results than just one.



GRU



Multilayer RNNs

In case of multiple layers, residual connections usually improve results. Because dimensionality has to be the same, they are usually applied from the second layer.



GRU



Bidirectional RNN

To consider both the left and right contexts, a *bidirectional* RNN can be used, which consists of parallel application of a *forward* RNN and a *backward* RNN.



The outputs of both directions can be either *added* or *concatenated*. Even if adding them does not seem very intuitive, it does not increase dimensionality and therefore allows residual connections to be used in case of multilayer bidirectional RNN.

Word Embeddings

Ú FA

We might represent *words* using one-hot encoding, considering all words to be independent of each other.

However, words are not independent – some are more similar than others.

Ideally, we would like some kind of similarity in the space of the word representations.

Distributed Representation

The idea behind distributed representation is that objects can be represented using a set of common underlying factors.

We therefore represent words as fixed-size *embeddings* into \mathbb{R}^d space, with the vector elements playing role of the common underlying factors.

These embeddings are initialized randomly and trained together with the rest of the network.

GRU

Word Embeddings

Ú_F^{AL}

The word embedding layer is in fact just a fully connected layer on top of one-hot encoding. However, it is not implemented in that way.

Instead, a so-called *embedding* layer is used, which is much more efficient. When a matrix is multiplied by an one-hot encoded vector (all but one zeros and exactly one 1), the row corresponding to that 1 is selected, so the embedding layer can be implemented only as a simple lookup.

In TensorFlow, the embedding layer is available as

tf.keras.layers.Embedding(input_dim, output_dim)

GRU

Word Embeddings

Even if the embedding layer is just a fully connected layer on top of one-hot encoding, it is important that this layer is *shared* across the whole network.



NPFL114, Lecture 7

GRU

HighwayNetworks

RNNRegularization

RNNApplications

WordEmbeddings

