

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Add WeChat [edu\\_assist\\_pro](#)  
Recurrent Neural Networks (RNN)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](#)

Long distance dependency in sequence labeling

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

Limitation of window based feature extraction for linear sequence models can reach a very high accuracy, but are insufficient cases:

- ▶ POS tagging

\_VBZ pianos

- ▶ Named Entity Recognition

<https://eduassistpro.github.io/>

"Pretty soon I don't know what but *I*-title is going to happen"

[Add WeChat edu\\_assist\\_pro](https://eduassistpro.github.io/)

- ▶ Language modeling: The man who whistles **tunes** pianos

- ▶ ...

## Simple Recurrent Neural Networks

- ▶ A simple recurrent neural network. It is defined as:

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

where  $\Theta \in \mathbb{R}^{(d_x+d_h) \times d_h}$ .  $g$  is a nonlinear transformation, usually hyperbolic tangent  $\tanh$ .

- ▶ Although each  $v_m$  only depends on  $h_{m-1}$ , but this vector is affected by all previous tokens  $w_0, w_1, \dots, w_{m-1}$ . This is crucially different from n-gram language models.
- ▶ While in principle this simple RNN can handle long distance dependencies, in practice it is quite inadequate, due to the repeated application of the non-linearity.

## Layers in an RNN

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ How many layers do we have in an RNN?  
<https://eduassistpro.github.io/>
- ▶ The number of layers in an RNN is not fixed and varies with the length of the sequence.
- ▶ Derivatives of the hidden states can be computed using backpropagation through time, which is implemented in various deep learning frameworks such as Torch, MXNet, and TensorFlow.
- ▶ How many weight matrices do we need to learn?

## Parameters of a simple RNN

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶  $\phi_i \in \mathbb{R}^K$ , the “input” word vectors (word embedding);
- ▶  $\beta_i \in \mathbb{R}^K$ , the “output” word vectors;
- ▶  $\Theta \in \mathbb{R}^{K \times K}$
- ▶  $h_0$ , the initial hidden state.

Each of these parameters can be tuned by formulating a loss function over the training corpus  $L(w)$

## Backpropagation through time (BPTT)

<https://eduassistpro.github.io/>

- ▶ Let  $\ell_{m+1}$  represent the negative log-likelihood of word  $m + 1$ ,

Add  $\frac{\partial \ell_{m+1}}{\partial \theta_{k,k'}} h_m$  to the gradient:

- ▶ Since the loss  $d$  (not through <https://eduassistpro.github.io/>)

Add  $\frac{\partial \ell_{m+1}}{\partial \theta_{k,k'}} h_m$  to the gradient:

- ▶ What is the derivative of this function with respect to an element in  $\Theta$ ?

## Recurrence in gradient computation

<https://eduassistpro.github.io/>

- ▶ Recall the simple RNN with the Elman Unit

**Assignment Project Exam Help**  
 $h_m = g(\Theta h_{m-1} + x_m)$

- ▶ For any element in the RNN, Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

<https://eduassistpro.github.io/>)

- ▶ Applying the chain rule and the product rule

Add WeChat [edu\\_assist\\_pro](https://edu_assist_pro)

$$\frac{\partial h_{m,k}}{\partial \theta_{k,k'}} = g'(x_{m,k} + \theta_k \cdot h_{m-1}) \left( h_{m-1,k'} + \theta_k \cdot \frac{\partial h_{m-1}}{\partial \theta_{k,k'}} \right)$$

- ▶ This means the derivative of  $\frac{\partial h_m}{\partial \theta_{k,k'}}$  depends on  $\frac{\partial h_{m-1}}{\partial \theta_{k,k'}}$ , which in turn depends on  $\frac{\partial h_{m-2}}{\partial \theta_{k,k'}}$ , etc., till  $h_0$  is reached

## Variants of RNNs

<https://eduassistpro.github.io/>

- ▶ An Elman RNN helps us demonstrate RNN's capacity of modeling longer context, and the needs to use the network [Assignment Project Exam Help](#) [Add WeChat edu\\_assist\\_pro](#)

<https://eduassistpro.github.io/>

- ▶ A more common RNN also parameterizes another parameter: [Assignment Project Exam Help](#) [Add WeChat edu\\_assist\\_pro](#)

$$\text{RNN}(\mathbf{x}_m, \mathbf{h}_{m-1}) = g(\Theta^h \mathbf{h}_{m-1} + \Theta^x \mathbf{x}_m)$$

Updating  $\Theta^x$  does not require backpropagate through time.

## Hyperparameters

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ The optimal size of the word and context vector  $K$  loosely depends on the training data.
- ▶ If the training data is large, a larger  $K$  is preferred.
- ▶ Conversely, a smaller  $K$  should be used, the model may memorize the training data at a faster rate.

## Vanishing or exploding gradients and gated networks

<https://eduassistpro.github.io/>

- ▶ RNNs require repeated applications of the non-linear functions. Backpropagation can lead to van (gradients decay to zero) or exploding gradients (increase towards infinity).
- ▶ Exploding gradients (set a maximum)
- ▶ Vanishing gradients must be addressed itself. Gated networks such as LSTM (Long short-term memory) and GRU are popular solutions to this problem.
- ▶ You can find a demonstration here: [https://cs224d.stanford.edu/notebooks/vanishing\\_grad\\_example.html](https://cs224d.stanford.edu/notebooks/vanishing_grad_example.html)

## Gated Recurrent Units (GRUs)

<https://eduassistpro.github.io/>

### Assignment Project Exam Help

- ▶ GRU uses a *Reset* gate and an *Update* to control how much information from the previous hidden state  $t$  next hidden state.

Reset gate

<https://eduassistpro.github.io/>

update gate

$$u_{m+1} = \sigma(\Theta^{h \rightarrow u} + b_u)$$

update candidate

$$\tilde{h}_{m+1} = \tanh(\Theta^h + \Theta^{x \rightarrow h} x_{m+1})$$

output

$$h_{m+1} = u_{m+1} \odot h_m + (1 - u_{m+1}) \tilde{h}_{m+1}$$

What is a gate?

<https://eduassistpro.github.io/>

## Assignment Project Exam Help

- ▶ “Hard” gate vs “soft” gate:

Add WeChat edu\_assist\_pro

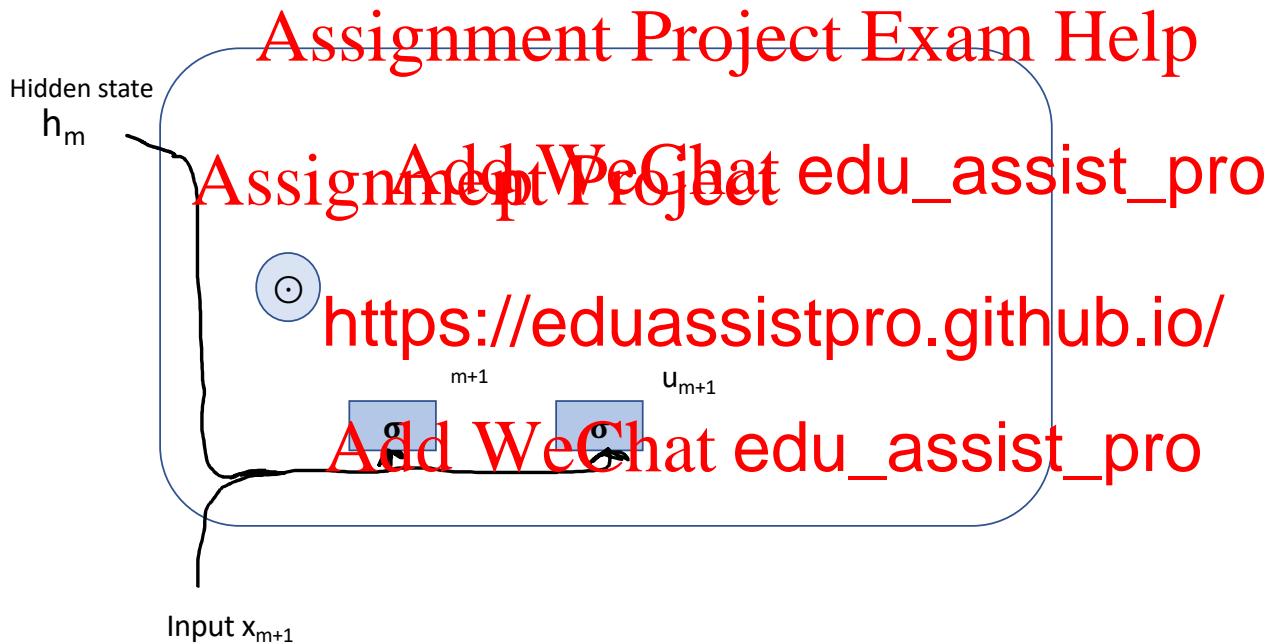
0  
0.998  
0  
0

Add WeChat edu\_assist\_pro

- ▶ Hard gates are easier to understand conceptually, but soft gates are differentiable (therefore learnable)

## GRU gates

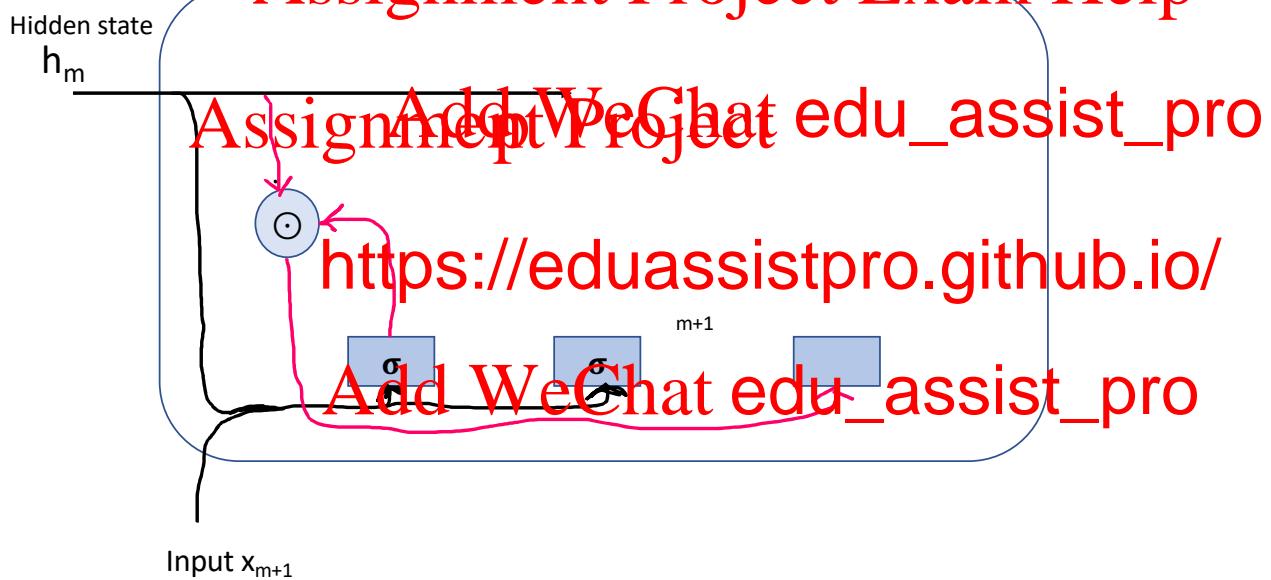
<https://eduassistpro.github.io/>



GRU gates

<https://eduassistpro.github.io/>

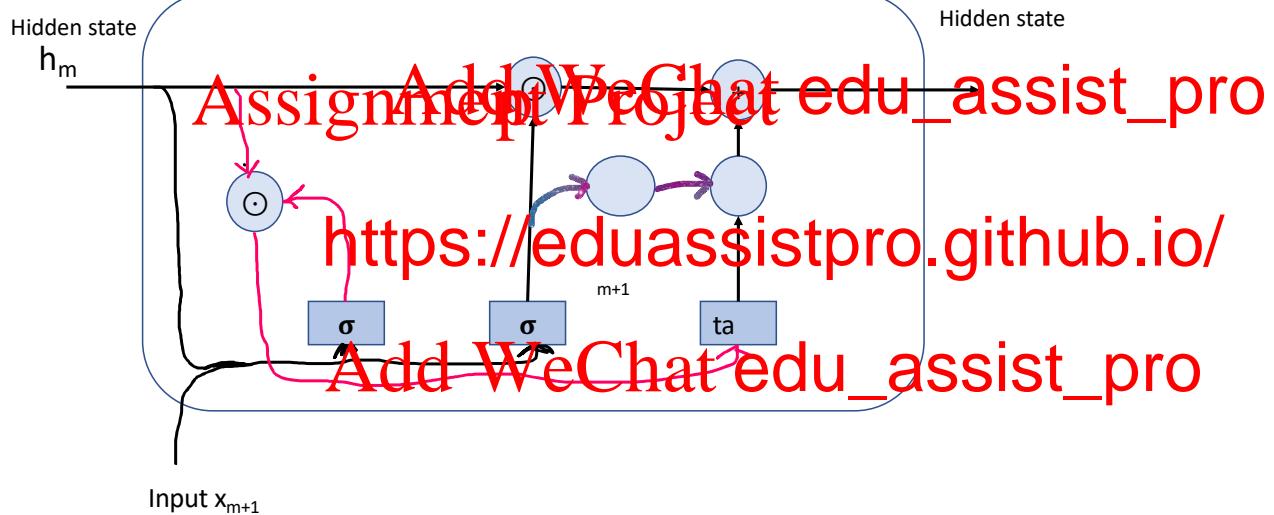
Assignment Project Exam Help



GRU gates

<https://eduassistpro.github.io/>

Assignment Project Exam Help



GRU: Pytorch implementation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

# Long Short Term Memory (LSTM)

<https://eduassistpro.github.io/>

- ▶ LSTM enhances the hidden  $\mathbf{h}_m$  with a memory cell  $\mathbf{c}_m$
- ▶ This memory cell does not pass through a squashing function ( $\tanh$  or  $\sigma$ ), and can propagate through the distances.

forget gate <https://eduassistpro.github.io/>

input gate  $i_{m+1} = \sigma(\Theta^h \mathbf{x}_{m+1} + \mathbf{b}_i)$

update candidate  $\tilde{\mathbf{c}}_{m+1} = \tanh(\Theta^x \mathbf{x}_{m+1})$

memory cell update  $\mathbf{c}_{m+1} = \mathbf{f}_{m+1} \odot \mathbf{c}_m + i_{m+1} \odot \tilde{\mathbf{c}}_{m+1}$

output gate  $\mathbf{o}_{m+1} = \sigma(\Theta^{h \rightarrow o} \mathbf{h}_m + \Theta^{x \rightarrow o} \mathbf{x}_{m+1} + \mathbf{b}_o)$

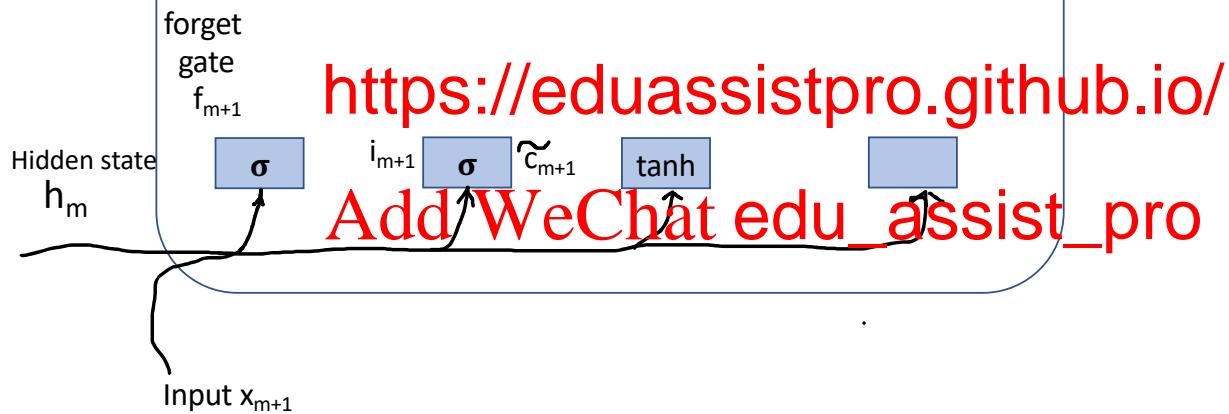
output  $\mathbf{h}_{m+1} = \mathbf{o}_{m+1} \odot \tanh(\mathbf{c}_{m+1})$

## LSTM gates

<https://eduassistpro.github.io/>

Assignment Project Exam Help

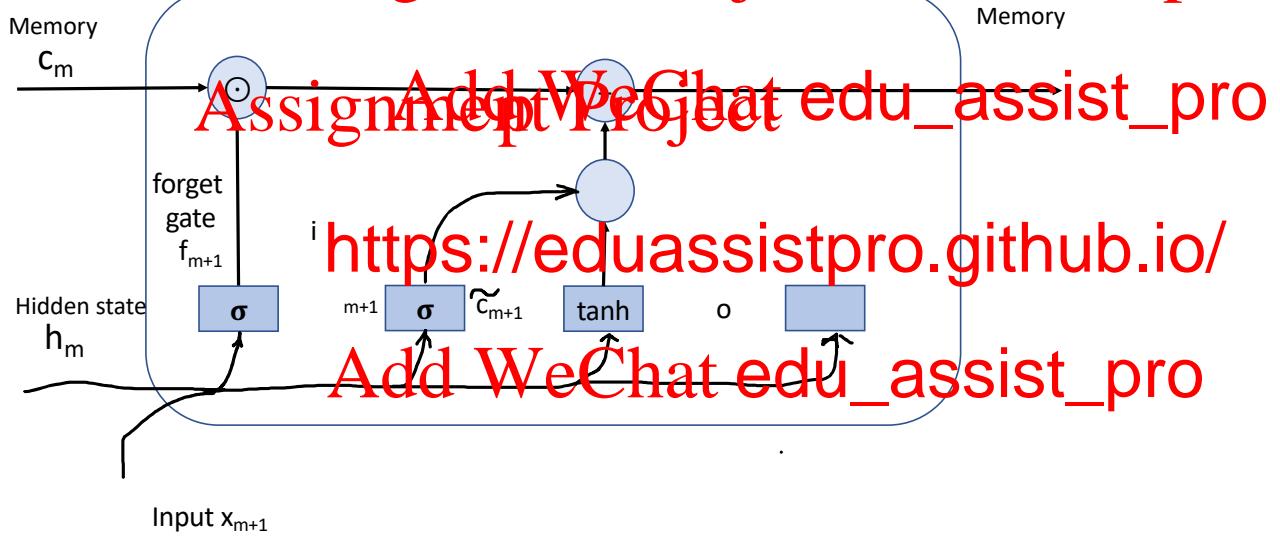
Assignment Project Exam Help Add WeChat edu\_assist\_pro



## LSTM gates

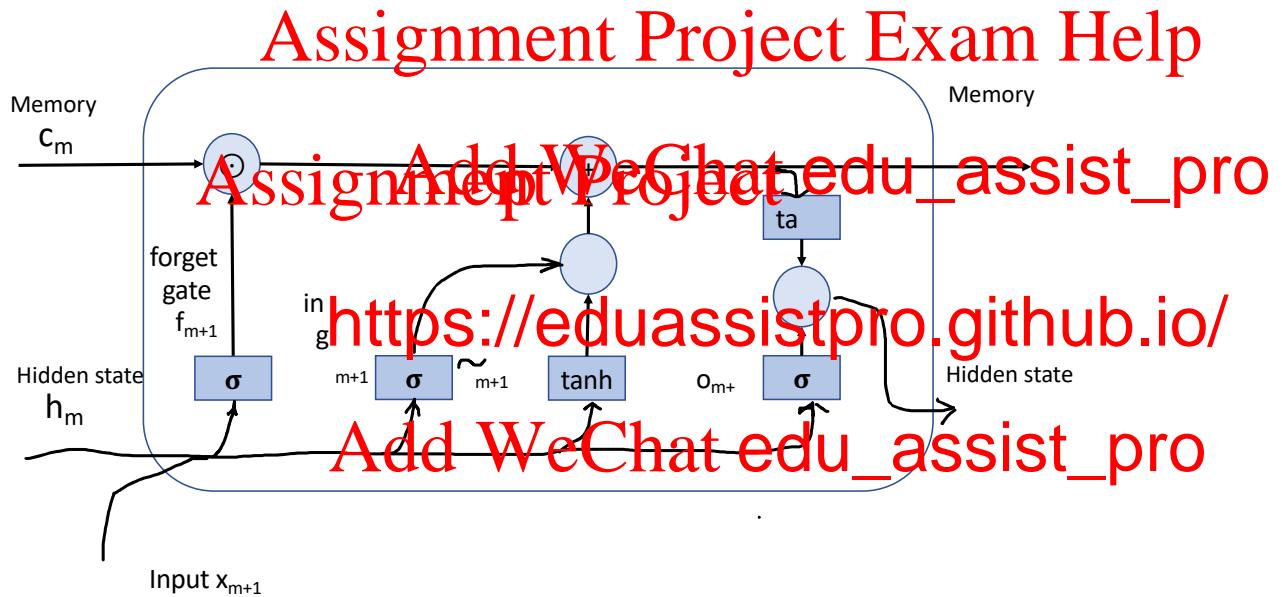
<https://eduassistpro.github.io/>

Assignment Project Exam Help



## LSTM gates

<https://eduassistpro.github.io/>



LSTM: Pytorch implementation

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Help WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

<https://eduassistpro.github.io/>

Add WeChat [edu\\_assist\\_pro](https://eduassistpro.github.io/)

The successful stories of Recurrent Neural networks in NLP

<https://eduassistpro.github.io/>

RNN is very versatile and can be used in a variety of different NLP problems:

- ▶ Sequence labeling problems: POS tagging, Recognitio
- ▶ Sequence to sequence problems: machine Translation, Machine Reading, Dialogue generation, etc.
- ▶ It also serves as an underlying feature extractor to many classification problems (e.g., the use of Bi-LSTM RNNs for extracting spans)

## RNNs for sequence labeling

- ▶ Recall RNN output at time step  $m$ :  
 $\text{https://eduassistpro.github.io/}$

Assignment Project Exam Help  
 $h_m = g(x_m, h_{m-1}), m = 1, 2, \dots, M$

- ▶ Scoring each label  $y_m$  as a linear function:  
 $\text{https://eduassistpro.github.io/}$

Assignment Project Exam Help  
Add WeChat edu\_assist\_pro

- ▶ We can also turn this score into a probability using SoftMax:  
 $\text{https://eduassistpro.github.io/}$

$$P(y|\mathbf{w}_{1:m}) = \frac{\exp \psi_m(y)}{\sum_{y' \in \mathcal{Y}} \exp \psi(y')}$$

- ▶ This is a classifier that uses only the context from the left.

## Bi-directional LSTM (BiLSTM)

<https://eduassistpro.github.io/>

- ▶ Forward LSTM and Backward LSTM

$\vec{h} = g(x_1, \vec{h}_{m-1}), m$   
 $\vec{h} = g(x_m, \vec{h}_{m+1}), m$

<https://eduassistpro.github.io/>

- ▶ BiLSTM summarizes the surrounding directions, typically better than windowing ngrams both sides of the target word.
- ▶ But it doesn't take into consideration the transitions between tags.

## Neural Structure Prediction: LSTM-CRF

Neural sequence labeling algorithm by defining the

Assignment Project Exam Help

$$s_m(y_m, y_{m-1}) = \beta \cdot h_m$$

Assignment Help WeChat edu\_assist\_pro

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

