

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Recurrent Neural Networks (RNNs) for sequence
labeling

<https://powcoder.com>

Add WeChat powcoder

Long distance dependency in sequence labeling

<https://powcoder.com>

Assignment Project Exam Help

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

- ▶ POS tagging: The man who whistles **tunes**_VBZ pianos
- ▶ Named Entity Recognition: Norma Jean's song entitled "Pretty soon I don't know what but **something**_I-title is going to happen"
- ▶ Language modeling: The man who whistles **tunes** pianos
- ▶ ...

Simple Recurrent Neural Networks

- ▶ A simple recurrent operation is called **Elman unit**. It is defined as:

$$\text{RNN}(x_m, \Theta) = g(\Theta h_{m-1} + x_m)$$

where $\Theta \in \mathbb{R}^{K \times K}$ is a recurrent matrix and g is a nonlinear transformation function, often defined as the element-wise hyperbolic tangent \tanh .

- ▶ Although each w_m only depends on the context vector 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://powcoder.com>

Assignment Project Exam Help

- ▶ How many layers do we have in an RNN?
 - ▶ The number of layers in an RNN is not fixed and varies with the length of the input.
 - ▶ Derivatives can be computed automatically with packages such as Torch, MXNet, and TensorFlow.
- ▶ How many weight matrices do we need for an RNN?

Parameters of a simple RNN

<https://powcoder.com>

Assignment Project Exam Help

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

Each of these parameters can be tuned by formulating an objective over the training corpus $L(\mathbf{w})$

Backpropagation through time (BPTT)

<https://powcoder.com>

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

$$\ell_{m+1} = -\log P(w_{m+1} | w_1, w_2, \dots, w_m)$$

- Since the loss depends on the parameters only through \mathbf{h}_m (not through x_m), we can apply the chain rule of differentiation:

$$\frac{\partial \ell_{m+1}}{\partial \theta_{k,k'}} = \frac{\partial \ell_{m+1}}{\partial \mathbf{h}_m} \frac{\partial \mathbf{h}_m}{\partial \theta_{k,k'}}$$

- What is the derivative of this function with respect to an element in Θ ?

Recurrence in gradient computation

<https://powcoder.com>

- Recall the simple RNN with the Elman Unit

Assignment Project Exam Help

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

- For any element $h_{m,k}$

Add WeChat powcoder

$$h_{m,k} = g(\theta_k \cdot \mathbf{h}_{m-1} + x_{m,k})$$

<https://powcoder.com>

- Applying the chain rule and the product rule, we get:

Add WeChat powcoder

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

- This means the derivative of $\frac{\partial \mathbf{h}_m}{\partial \theta_{k,k'}}$ depends on $\frac{\partial \mathbf{h}_{m-1}}{\partial \theta_{k,k'}}$, which in turn depends on $\frac{\partial \mathbf{h}_{m-2}}{\partial \theta_{k,k'}}$, etc., till \mathbf{h}_0 is reached

Variants of RNNs

<https://powcoder.com>

- ▶ An Elman RNN helps us demonstrate RNNs capacity of modeling longer context, and the needs to use BPTT update the network.

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

- ▶ A more common RNN also parameterize the input with another parameter:

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

Updating Θ^x does not require backpropate through time.

Hyperparameters

<https://powcoder.com>

Assignment Project Exam Help

- ▶ The optimal size of the word and context vectors K loosely depends on the size of the training data.
- ▶ If the training data is large, a larger K is preferred.
- ▶ Conversely, a smaller K should be chosen. Otherwise, the model may memorize the training data and doesn't generalize.

Vanishing or exploding gradients and gated networks

<https://powcoder.com>

- ▶ RNNs require repeated applications of the non-linear functions. Backpropagation can lead to vanishing gradients (gradients decay to zero) or exploding gradients (gradients increase towards infinity).
- ▶ Exploding gradients can be addressed by clipping gradients (set a maximum value)
- ▶ Vanishing gradients must be addressed by changing the model itself. Gated networks such as LSTM (Long Term 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

Gated Recurrent Units (GRUs)

<https://powcoder.com>

Assignment Project Exam Help

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

Add WeChat powcoder

Reset gate

$$r_{m+1} = \sigma(\Theta^{h \rightarrow r} \mathbf{h}_m + \Theta^{x \rightarrow r} \mathbf{x}_{m+1} + \mathbf{b}_r)$$

update gate

$$u_{m+1} = \sigma(\Theta^{h \rightarrow u} \mathbf{h}_m + \Theta^{x \rightarrow u} \mathbf{x}_{m+1} + \mathbf{b}_u)$$

update candidate

$$\tilde{\mathbf{h}}_{m+1} = \tanh(\Theta^{h \rightarrow h} (\mathbf{h}_m \odot \mathbf{r}_{m+1}) + \Theta^{x \rightarrow h} \mathbf{x}_{m+1})$$

output

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

What is a gate?

<https://powcoder.com>

Assignment Project Exam Help

- “Hard” gate vs “soft” gate:

Assignment Project Exam Help

$$\delta(x) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \quad \sigma(x) = \begin{bmatrix} 0 \\ 0.998 \\ 0 \\ 0 \end{bmatrix}$$

Add WeChat powcoder

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

GRU gates

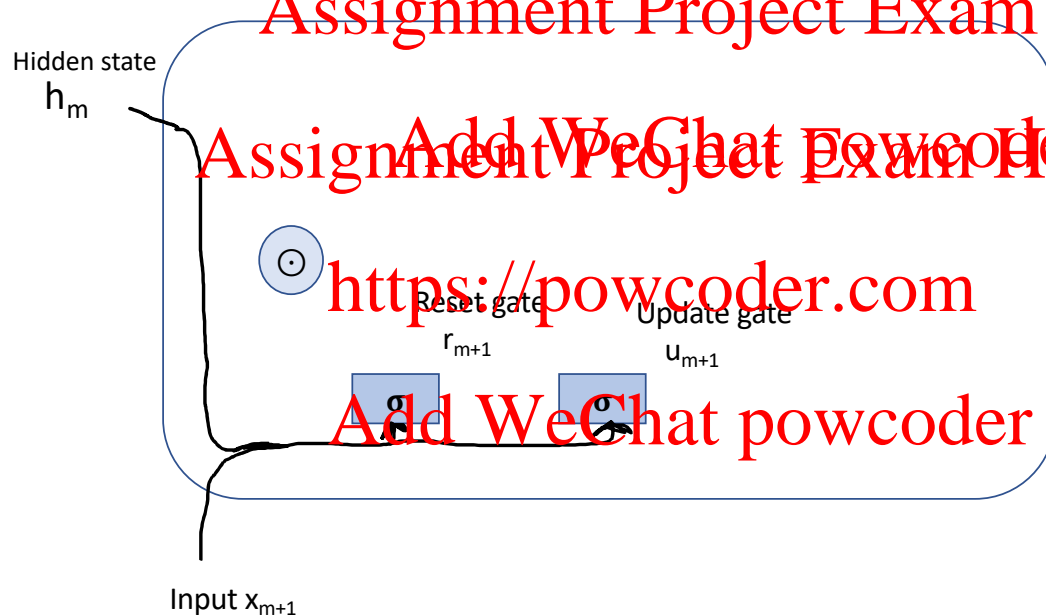
<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



GRU gates

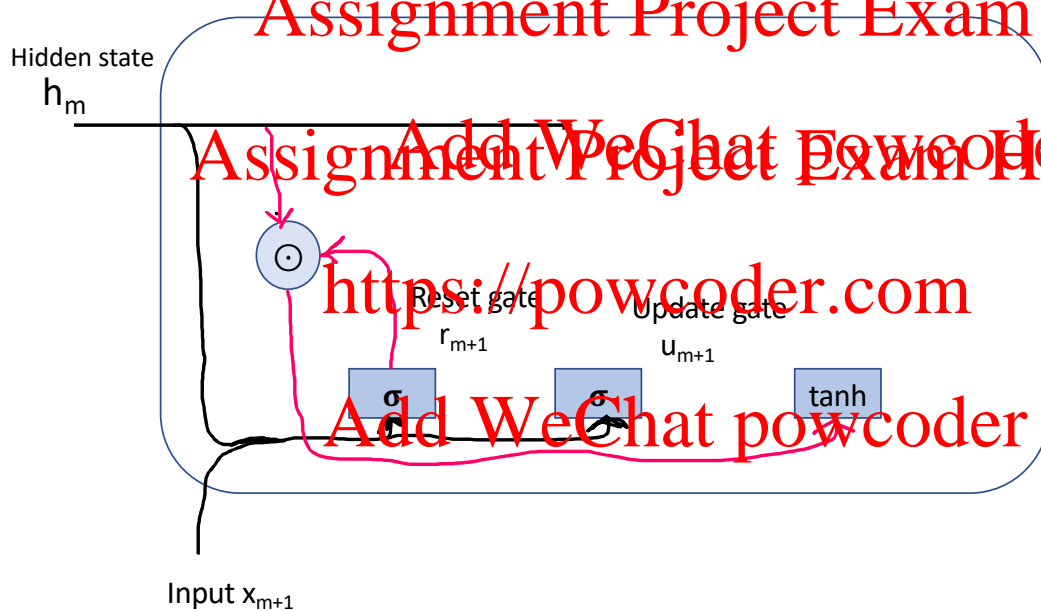
<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

<https://powcoder.com>

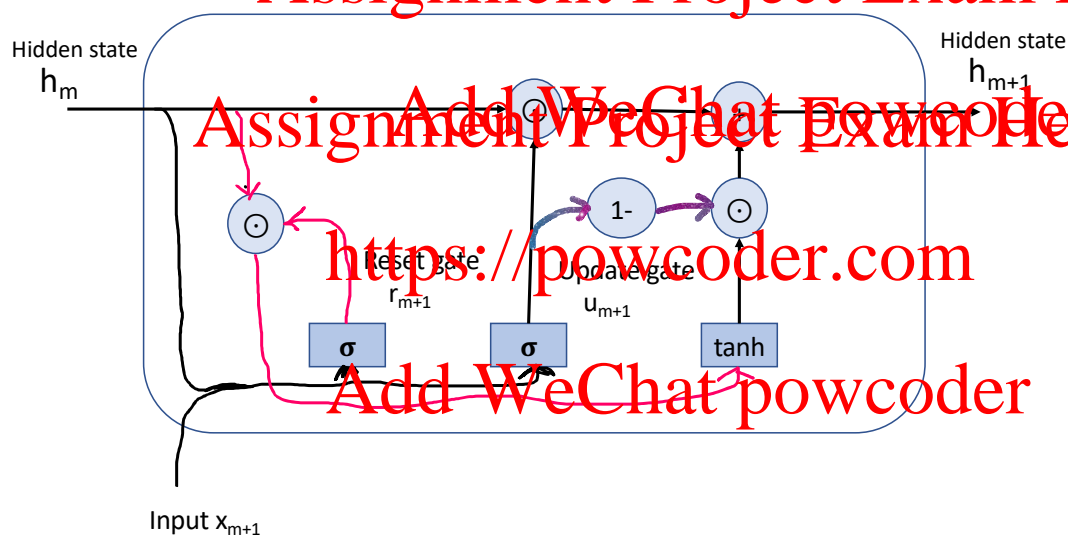
Add WeChat powcoder



GRU gates

<https://powcoder.com>

Assignment Project Exam Help



GRU: Pytorch implementation

<https://powcoder.com>

nn.GRU Assignment Project Exam Help

```
In [31]: gru = nn.GRU(3, 6, 2) #input dimension is 3, hidden dimension is 6, and 2 layers
pretrained_weights = torch.FloatTensor([[1,3,5],[2,6,7],[5,8,1,9],[4, 3, 7]])
embed_pretrained = nn.Embedding.from_pretrained(pretrained_weights)
input = torch.FloatTensor([1,2],[0,1],[0,5,7,2,1]) #sequence length = batch size: 2,
input_embedded = embed_pretrained(input) #embedding size: 3
h0 = torch.randn(2, 2, 6)
gru_output, hn = gru(input_embedded,h0)
print(gru_output)

tensor([[[[-0.5846, -0.5552,  0.0911, -0.7767,  0.7791,  0.5191],
          [ 0.4765,  0.0570,  0.6168,  0.8534, -0.3751,  0.0189]],

         [[-0.4376, -0.5312, -0.3698, -0.5889,  0.3567,  0.3723],
          [ 0.278, -0.1900,  0.812,  0.4208, -0.3358,  0.2713]],

         [[-0.3307, -0.5805, -0.5468, -0.5139,  0.1180,  0.3092],
          [ 0.1376, -0.3102, -0.2409,  0.1540, -0.6041,  0.3456]],

         [[-0.2772, -0.5828, -0.6242, -0.5065, -0.0066,  0.2383],
          [ 0.0673, -0.3882, -0.4041, -0.0143, -0.6205,  0.3542]]],
        grad_fn=<StackBackward>)
```

<https://powcoder.com>

Add WeChat powcoder

Long Short Term Memory (LSTM)

<https://powcoder.com>

- ▶ 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 σ), and can propagate through the network over long distances.

forget gate $\mathbf{f}_{m+1} = \sigma(\Theta^{h \rightarrow f} \mathbf{h}_m + \Theta^{x \rightarrow f} \mathbf{x}_{m+1} + \mathbf{b}_f)$

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

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

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

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

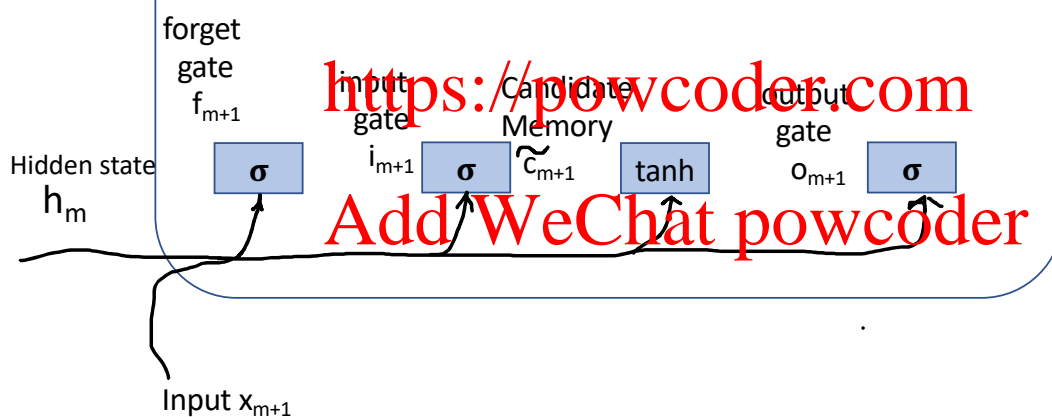
<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

<https://powcoder.com>

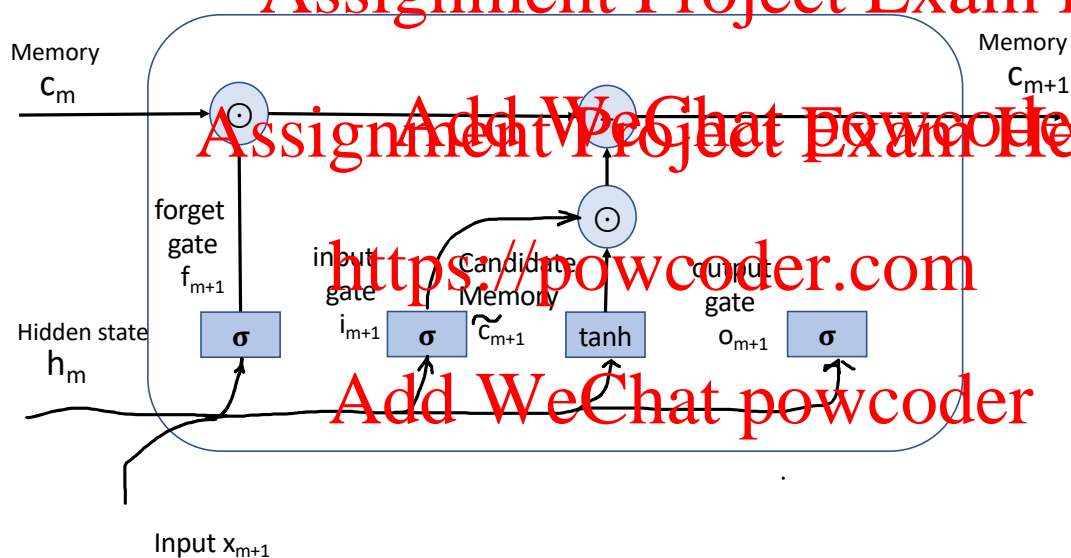
Add WeChat powcoder



LSTM gates

<https://powcoder.com>

Assignment Project Exam Help



Assignment Project Exam Help

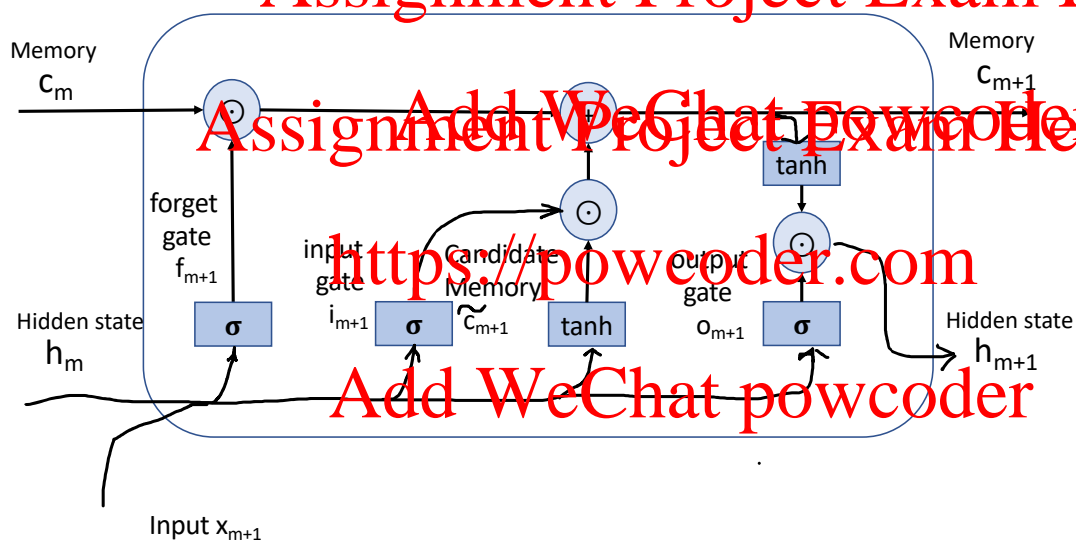
<https://powcoder.com>

Add WeChat powcoder

LSTM gates

<https://powcoder.com>

Assignment Project Exam Help



LSTM: Pytorch implementation

<https://powcoder.com>

nn.LSTM

- nn.LSTM parameters: input feature size *input_size*, hidden layer size *hidden_size*, layers *num_layers*

```
In [30]: lstm = nn.LSTM(3, 6, 2) #input dimension is 3, hidden dimension is 6, and 2 layers
pretrained_weights = torch.FloatTensor([[[1,3,5],[2,6,7],[5,8,1,9],[4, 3, 7]])
embed_pretrained = nn.Embedding.from_pretrained(pretrained_weights)
input = torch.LongTensor([[[[1,2,3],[4,5,6],[7,8,9]]]]) #seq: [1,2,3,4,5,6,7,8,9], batch size: 2,
#Dimensions of input batch: sequence length, batch size, and input feature size
input_embedded = embed_pretrained(input) #embedding size: 3
#(initialization) hidden state and cell state dimensions:
#number of layers times num. of directions l x d, batch size b*, hidden dimensions *h
h0 = torch.randn(2, 2, 6)
c0 = torch.randn(2, 2, 6)
#print(h0)
lstm_output, (hn,cn) = lstm(input_embedded, (h0,c0))
print(lstm_output)

tensor([[[[-0.1751, -0.0039,  0.0705,  0.1108, -0.1001,  0.2951],
          [-0.0536, -0.1103,  0.1305, -0.0536, -0.0359, -0.2782]],

          [[-0.1298, -0.0287,  0.0813,  0.0542,  0.1038, -0.1680],
          [-0.0892, -0.0794,  0.1390,  0.0033,  0.1406, -0.2525]],

          [[-0.1128, -0.0194,  0.0992,  0.0692,  0.1854, -0.1299],
          [-0.0747, -0.0993,  0.1663,  0.0730,  0.2503, -0.1682]],

          [[-0.0880, -0.0654,  0.1431,  0.0854,  0.2612, -0.1044],
          [-0.0714, -0.0684,  0.1699,  0.0885,  0.2691, -0.1043]]],

        grad_fn=<StackBackward>)
```

The successful stories of Recurrent Neural networks in NLP

<https://powcoder.com>

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

- ▶ Sequence labeling problems: POS tagging, Named Entity Recognition, etc.
- ▶ Sequence to sequence translation problems: Neural Machine Translation, Machine Reading, Dialogue systems, language generation, etc.
- ▶ It also serves as an underlying feature extraction mechanism to many classification problems (e.g., the use of Bi-LSTM RNNs for extracting spans)

RNNs for sequence labeling

- Recall RNN outputs a hidden vector \mathbf{h}_m at each time step:

$$\mathbf{h}_m = g(\mathbf{x}_m, \mathbf{h}_{m-1}), m = 1, 2, \dots, M$$

- Scoring each label y_m as a linear function of \mathbf{h}_m

$$\psi_m(y) = \beta_y \cdot \mathbf{h}_m$$
$$\hat{y} = \underset{y}{\operatorname{argmax}} \psi_m(y)$$

- We can also turn this score into a probabilistic distribution using SoftMax:

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

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

Bi-directional LSTM (BiLSTM)

<https://powcoder.com>

- Forward LSTM and Backward LSTM

$$\begin{aligned}\vec{h}_m &= g(\mathbf{x}_m, \vec{h}_{m-1}), m = 1, 2, \dots, M \\ \overleftarrow{h}_m &= g(\mathbf{x}_m, \overleftarrow{h}_{m+1}), m = 1, 2, \dots, M\end{aligned}$$

<https://powcoder.com>

- BiLSTM summarizes the surrounding context from both directions, typically better than window-based context using ngrams both sides of the target word.
- But it doesn't taking into consideration the transitions between tags.

Neural Structure Prediction: LSTM-CRF

Neural sequence labeling can be combined with the Viterbi algorithm by defining the local score as:

$$s_m(y_m, y_{m-1}) = \beta \cdot \mathbf{h}_m + \eta_{y_{m-1}, y_m}$$

