

CS 4774 Machine Learning

CNNs and RNNs

Yangfeng Ji

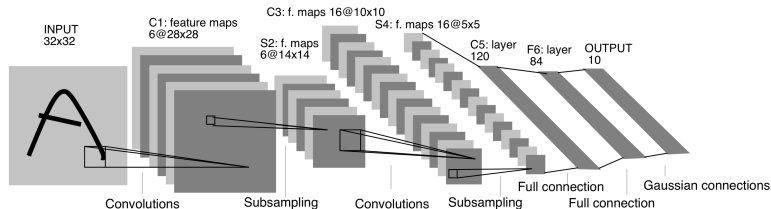
Information and Language Processing Lab
Department of Computer Science
University of Virginia



1. Convolutional Neural Networks
2. Recurrent Neural Networks
3. RNN Language Modeling
4. Challenge of Training RNNs

Convolutional Neural Networks

A classical neural network architecture designed for handwritten and machine-printed character recognition.



This architecture repeats the two components twice before connecting with a fully-connected layer

- ▶ convolutional layer
- ▶ subsampling (pooling) layer

[LeCun et al., 1998]

Convolutional Operations

1-D convolutional operation is defined as

$$c_j = \mathbf{m}^\top \mathbf{x}_{j:j+n-1} \quad (1)$$

where $\mathbf{m} \in \mathbb{R}^n$ is convolutional filter with window size n , $\mathbf{x} \in \mathbb{R}^T$ input signal with size T , and $j \leq T - n + 1$.

Convolutional Operations

1-D convolutional operation is defined as

$$c_j = \mathbf{m}^T \mathbf{x}_{j:j+n-1} \quad (1)$$

where $\mathbf{m} \in \mathbb{R}^n$ is convolutional filter with window size n , $\mathbf{x} \in \mathbb{R}^T$ input signal with size T , and $j \leq T - n + 1$.

Example

With

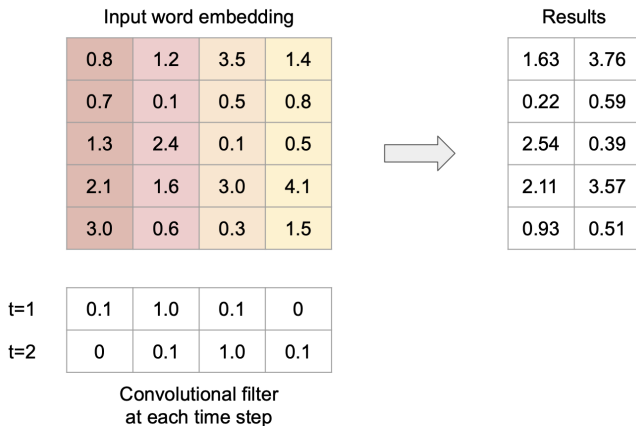
- ▶ $\mathbf{m} = [m_1, m_2, m_3]^T$, and
- ▶ $\mathbf{x} = [x_1, x_2, x_3, x_4, \dots, x_T]^T$,

when $j = 2$

$$c_2 = m_1 x_2 + m_2 x_3 + m_3 x_4 \quad (2)$$

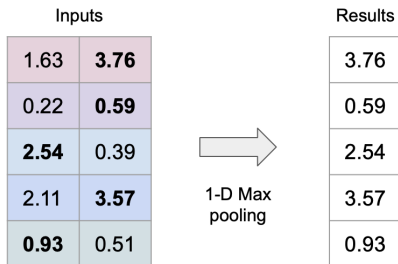
Convolutional Operations (II)

An example of 1-D convolutional operations



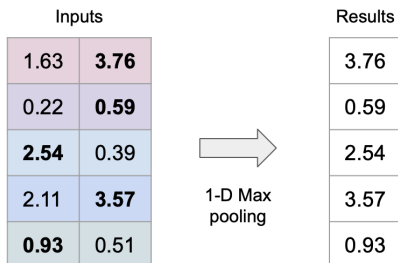
There are three popularly used pooling techniques

- ▶ Max pooling



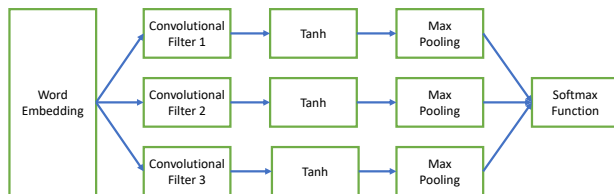
There are three popularly used pooling techniques

- ▶ Max pooling



- ▶ Average pooling [LeCun et al., 1998]
- ▶ Min pooling

A simple and effective convolutional neural network architecture for text classification [Kim, 2014]



- ▶ `torch.nn.Conv1d`: convolutional operation on each dimension of the word embeddings (no cross-dimension convolution)
- ▶ `torch.tanh`
- ▶ `torch.max`: max pooling on each dimension of the word embeddings
- ▶ `torch.cat`: concatenate three vectors from max pooling to form one single vector
- ▶ In actual implementation, the input is a 3-D tensor instead of a 2-D matrix

Advantages of CNNs

Comparing to Feed-forward NNs: Parameter sharing, sparse connections

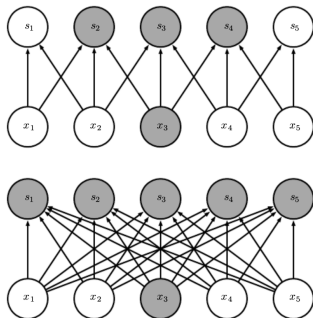


Figure: (1) upper plot: convolutional layer; (2) lower plot: fully-connected layer.

[Goodfellow et al., 2016]

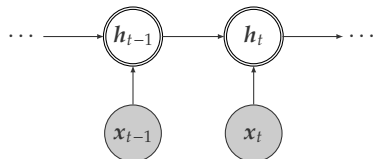
Recurrent Neural Networks

Recurrent Neural Networks (RNNs)

A simple RNN is defined by the following recursive function

$$h_t = f(x_t, h_{t-1}) \quad (3)$$

and depicted as



where

- ▶ h_{t-1} : hidden state at time step $t - 1$
- ▶ x_t : input at time step t
- ▶ h_t : hidden state at time step t

A Simple Transition Function

In the simplest case, the transition function f is defined with an **element-wise** Sigmoid function and a linear transformation of x_t and h_{t-1}

$$h_t = f(x_t, h_{t-1}) = \sigma(\mathbf{W}_h h_{t-1} + \mathbf{W}_i x_t + \mathbf{b}) \quad (4)$$

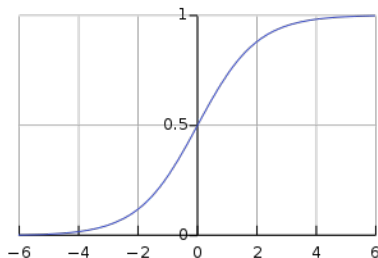
where

- ▶ x_t : input word embedding
- ▶ h_{t-1} : hidden state from previous time step
- ▶ \mathbf{W}_h : parameter matrix for **hidden states**
- ▶ \mathbf{W}_i : parameter matrix for **inputs**
- ▶ \mathbf{b} : bias term (also a parameter)

Sigmoid Function

A Sigmoid function with one-dimensional input $x \in (-\infty, \infty)$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



The potential numeric issue caused by the Sigmoid function

- ▶ $\sigma(x) \rightarrow 1$ with $x \gg 6$
- ▶ $\sigma(x) \rightarrow 0$, $x \ll -6$

The output of the Sigmoid function will approximate a constant, when the input value is beyond certain ranges

Unfolding RNNs

We can unfold this recursive definition of a RNN

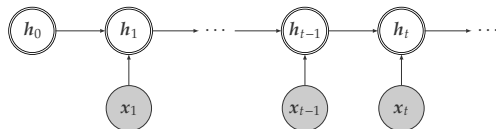
$$h_t = f(x_t, h_{t-1}) \quad (5)$$

Unfolding RNNs

We can unfold this recursive definition of a RNN

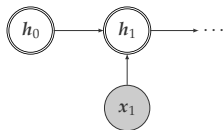
$$h_t = f(x_t, h_{t-1}) \quad (5)$$

as



$$\begin{aligned} h_t &= f(x_t, f(x_{t-1}, h_{t-2})) \\ &= f(x_t, f(x_{t-1}, f(x_{t-2}, h_{t-3}))) \\ &= \dots \\ &= f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(x_1, h_0) \dots))) \end{aligned} \quad (6)$$

Base condition defines the starting point of the recursive computation

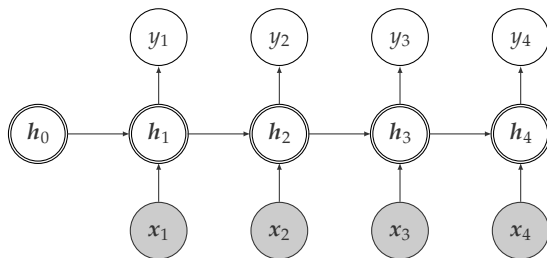


$$h_t = f(x_t, f(x_{t-1}, f(x_{t-2}, \dots f(x_1, h_0) \dots))) \quad (7)$$

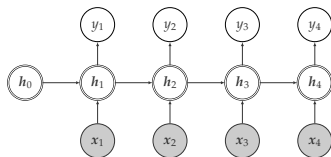
- ▶ h_0 : zero vector or parameter
- ▶ x_1 : input at time $t = 1$

RNN for Sequential Prediction

In general, RNNs can be used for any sequential modeling tasks



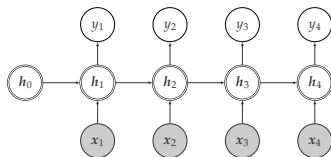
Sequential Modeling as Classification



- Prediction at each time step t

$$\hat{y}_t = \operatorname{argmax}_y P(y; \mathbf{h}_t) \quad (8)$$

Sequential Modeling as Classification



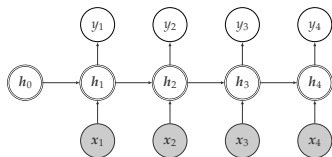
- Prediction at each time step t

$$\hat{y}_t = \operatorname{argmax}_y P(y; \mathbf{h}_t) \quad (8)$$

- Loss at single time step t

$$L_t(y_t, \hat{y}_t) = -\log P(y_t; \mathbf{h}_t) \quad (9)$$

Sequential Modeling as Classification



- ▶ Prediction at each time step t

$$\hat{y}_t = \operatorname{argmax}_y P(y; \mathbf{h}_t) \quad (8)$$

- ▶ Loss at single time step t

$$L_t(y_t, \hat{y}_t) = -\log P(y_t; \mathbf{h}_t) \quad (9)$$

- ▶ The total loss

$$\ell = \sum_{t=1}^T L_t(y_t, \hat{y}_t) \quad (10)$$

RNN Language Modeling

A language model defines the probability of x_t given $\mathbf{x} = (x_1, x_2, \dots, x_{t-1})$ as

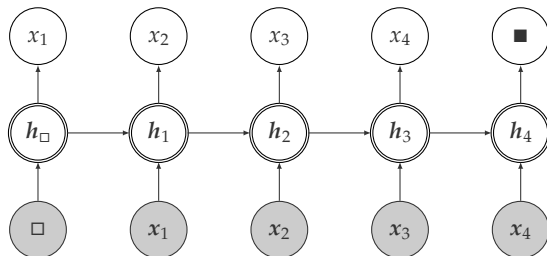
$$P(x_t \mid x_1, \dots, x_{t-1}) \tag{11}$$

and the joint probability as

$$\begin{aligned} P(\mathbf{x}_{1:T}) &= P(x_1) \cdot P(x_2 \mid x_1) \\ &\quad \cdot \dots \cdot \\ &\quad \cdot P(x_T \mid x_1, x_2, \dots, x_{T-1}) \end{aligned}$$

Language Modeling with RNNs

Using RNNs for language modeling

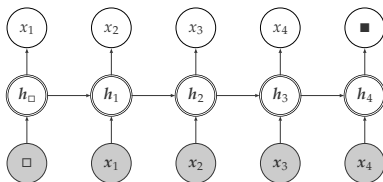


with two special tokens

$$\{\square, x_1, \dots, x_T, \blacksquare\}$$

RNN Language Models

For a given sentence $\{x_1, \dots, x_t\}$, the input at time t is **word embedding** x_t



The probability distribution of next word X_t

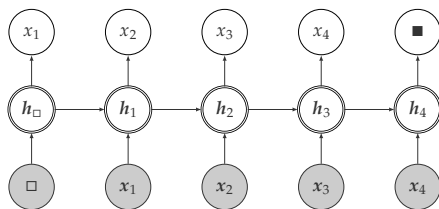
$$P(X_t = x \mid \mathbf{x}_{1:t-1}) = \frac{\exp(\mathbf{w}_{o,x}^\top \mathbf{h}_{t-1})}{\sum_{x' \in \mathcal{V}} \exp(\mathbf{w}_{o,x'}^\top \mathbf{h}_{t-1})} \quad (12)$$

where

- ▶ $\mathbf{w}_{o,x}$ is the output weight vector (parameter) associated with word x
- ▶ \mathcal{V} is the word vocabulary

Special Cases

Similar to statistical language modeling, there are also two special cases that we need to consider



$$\{\square, x_1, \dots, x_T, \blacksquare\}$$

The corresponding prediction functions are defined as

- ▶ At time $t = 1$

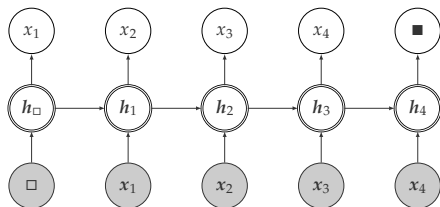
$$P(X_1 = x) \propto \exp(\mathbf{w}_{o,x}^\top \mathbf{h}_\square) \quad (13)$$

- ▶ At time $t = T$

$$P(X_T = \blacksquare \mid x_{1:T-1}) \propto \exp(\mathbf{w}_{o,x}^\top \mathbf{h}_{T-1}) \quad (14)$$

Challenge of Training RNNs

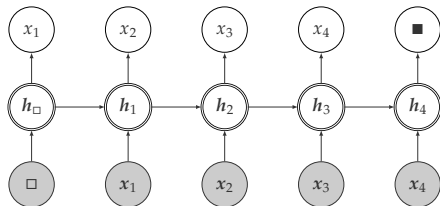
The training objective for each timestep is to predict the next token in the text



- ▶ Prediction at step t , $P(X_t = x \mid \mathbf{x}_{1:t-1}) = \frac{\exp(\mathbf{w}_{0,x}^T \mathbf{h}_{t-1})}{\sum_{x' \in \mathcal{V}} \exp(\mathbf{w}_{0,x'}^T \mathbf{h}_{t-1})}$
- ▶ Loss at step t , $L_t = -\log P(X_t = x \mid \mathbf{x}_{1:t-1})$

Let θ denote **all** model parameters

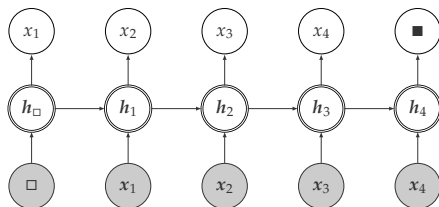
$$\frac{\partial \ell}{\partial \theta} = \sum_{t=1}^T \frac{\partial L_t}{\partial \theta} \quad (15)$$



Backpropagation Through Time [Rumelhart et al., 1985, BPTT]

Model Parameters

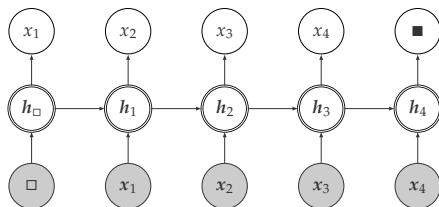
Before computing the gradient of each L_t with respect to model parameters, let us count how many parameters that we need consider



- ▶ Output parameter matrix $W_o = (w_{o,1}, \dots, w_{o,V})$

Model Parameters

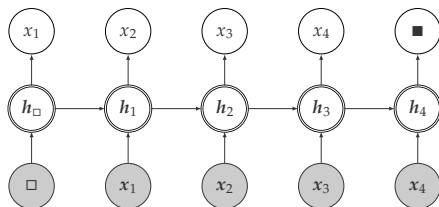
Before computing the gradient of each L_t with respect to model parameters, let us count how many parameters that we need consider



- ▶ Output parameter matrix $W_o = (w_{o,1}, \dots, w_{o,V})$
- ▶ Input word embedding matrix $X = (x_1, \dots, x_V)$

Model Parameters

Before computing the gradient of each L_t with respect to model parameters, let us count how many parameters that we need consider



- ▶ Output parameter matrix $W_o = (w_{o,1}, \dots, w_{o,V})$
- ▶ Input word embedding matrix $X = (x_1, \dots, x_V)$
- ▶ Neural network parameters W_h, W_i, b

Backpropagation Through Time

Take time step t as an example, we can take a look the gradient computation of some specific parameters

- ▶ Output model parameter $\frac{\partial L_t}{\partial w_{o,r}}$.

Backpropagation Through Time

Take time step t as an example, we can take a look the gradient computation of some specific parameters

- ▶ Output model parameter $\frac{\partial L_t}{\partial w_o}$.
- ▶ Neural network parameters, for example \mathbf{W}_h

$$\frac{\partial L_t}{\partial \mathbf{W}_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial \mathbf{h}_t} \cdot \left(\prod_{j=i}^{t-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_h} \right\} \quad (16)$$

Similar patterns for the other two neural network parameters \mathbf{W}_i and \mathbf{b}

Backpropagation Through Time

Take time step t as an example, we can take a look the gradient computation of some specific parameters

- ▶ Output model parameter $\frac{\partial L_t}{\partial w_o}$.
- ▶ Neural network parameters, for example \mathbf{W}_h

$$\frac{\partial L_t}{\partial \mathbf{W}_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial \mathbf{h}_t} \cdot \left(\prod_{j=i}^{t-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_h} \right\} \quad (16)$$

Similar patterns for the other two neural network parameters \mathbf{W}_i and \mathbf{b}

- ▶ Word embedding $\frac{\partial L_t}{\partial \mathbf{x}_{t'}}$
 - ▶ E.g., word embedding $\mathbf{x}_{t'}$ is the input of \mathbf{h}_t if $t' \leq t$, so ...

For each timestep, we need to compute the gradient using the chain rule:

$$\frac{\partial L_t}{\partial \mathbf{W}_h} = \sum_{i=1}^t \left\{ \frac{\partial L_t}{\partial \mathbf{h}_t} \cdot \left(\prod_{j=i}^{t-1} \frac{\partial \mathbf{h}_{j+1}}{\partial \mathbf{h}_j} \right) \cdot \frac{\partial \mathbf{h}_i}{\partial \mathbf{W}_h} \right\} \quad (17)$$

The chain rule of gradient will cause two potential problems in training RNNs

- ▶ vanishing gradients: $\frac{\partial L_t}{\partial \theta} \rightarrow 0$
- ▶ exploding gradients: $\frac{\partial L_t}{\partial \theta} \geq M$

[Pascanu et al., 2013]

Solution: **norm clipping** [Pascanu et al., 2013].

Consider the gradient $\mathbf{g} = \frac{\partial \ell}{\partial \theta}$,

$$\hat{\mathbf{g}} \leftarrow \tau \cdot \frac{\mathbf{g}}{\|\mathbf{g}\|} \quad (18)$$

when $\|\mathbf{g}\| > \tau$.

- ▶ Usually, $\tau = 3$ or 5 in practice.
- ▶ Smaller gradient will cause slower learning progress

Solution:

- ▶ initialize parameters carefully
- ▶ replace hidden state transition function $\sigma(\cdot)$ with other options

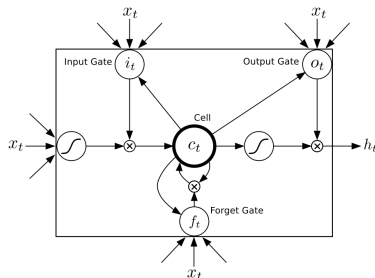
$$f(\mathbf{x}_t, \mathbf{h}_{t-1}) = \sigma(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_i \mathbf{x}_t + \mathbf{b}) \quad (19)$$

- ▶ LSTM [Hochreiter and Schmidhuber, 1997]
- ▶ GRU [Cho et al., 2014]

Long Short-Term Memory

Rather than directly taking input and hidden state as simple transition function, LSTM relies on three gates to control *how much* information it should take from input and hidden state before combining them together

$$\begin{aligned}i_t &= \sigma(\mathbf{W}_{xi}x_t + \mathbf{W}_{hi}h_{t-1} + \mathbf{W}_{ci}c_{t-1} + \mathbf{b}_i) \\f_t &= \sigma(\mathbf{W}_{xf}x_t + \mathbf{W}_{hf}h_{t-1} + \mathbf{W}_{cf}c_{t-1} + \mathbf{b}_f) \\c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(\mathbf{W}_{xc}x_t + \mathbf{W}_{hc}h_{t-1} + \mathbf{b}_c) \\o_t &= \sigma(\mathbf{W}_{xo}x_t + \mathbf{W}_{ho}h_{t-1} + \mathbf{W}_{co}c_t + \mathbf{b}_o) \\h_t &= o_t \circ \tanh(c_t)\end{aligned}$$



where \circ is the element-wise multiplication, $\{\mathbf{W}\}$ and $\{\mathbf{b}\}$ are parameters.

[Graves, 2013]



Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014).
On the properties of neural machine translation: Encoder-decoder approaches.
arXiv preprint arXiv:1409.1259.



Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016).
Deep Learning, volume 1.
MIT press Cambridge.



Graves, A. (2013).
Generating sequences with recurrent neural networks.
arXiv preprint arXiv:1308.0850.



Hochreiter, S. and Schmidhuber, J. (1997).
Long short-term memory.
Neural computation, 9(8):1735–1780.



LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. (1998).
Gradient-based learning applied to document recognition.
Proceedings of the IEEE, 86(11):2278–2324.



Pascanu, R., Mikolov, T., and Bengio, Y. (2013).
On the difficulty of training recurrent neural networks.
In Dasgupta, S. and McAllester, D., editors, *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1310–1318, Atlanta, Georgia, USA. PMLR.



Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1985).
Learning internal representations by error propagation.
Technical report, California Univ San Diego La Jolla Inst for Cognitive Science.