

Introduction to Deep Learning

Recurrent Neural Networks



maxwellcai.com



Wouldn't it be easier to handle sequences (1D) than images (2D)?



Yes and no...

Sequence Modeling

Language translation Speech recognition Music generation Sentiment classification Video captioning Chatbot

Text summary





Sequential Data

Sequential data are represented with 1D or 2D arrays.



Time	24/3	25/3	26/3	27/3	28/3	
Value	19.722	21.050	21.188	22.532	21.678	



In natural language processing (NLP),

- Each word is converted into a token;
- A token is typically represented as a numerical value;
- A sentence is a sequence of tokens;
- A sequence can have arbitrary length;
- Tokens may be **logically related** to each other.





Sequential Data Modeling

Naive Approaches



- Modeling the **full sentence** using a fully connected network requires a huge number of parameters!
- A FC network needs to model the pairwise relation between each word, thus the network scales as $\mathcal{O}(k^n)$, where k is size of the vocabulary and n is the length of the sequence.
- For example, consider that in the English language, there are k = 2000 words. To model a sentence of n = 100 words, one single FC layer requires 2000^{100} parameters!!



Sequence modeling with 1D convolution

CNNs works because they can take into account the relation between **nearby** tokens.



Image source: "Deep Learning With Python" (Francois Chollet)



Example: human activity recognition



The DNN can be trained with **existing** open datasets! e.g., https://archive.ics.uci.edu/ml/datasets/human+activity+recognition+using+smartphones

Human activity recognition with 1D Convolution

Layer 4 Convolutional layer

Layer 5 Convolutional layer

Layer 6 Dense layer

Layer 7 Dense layer

Ordóñez et al. (2016)

Sequence modeling with 1D convolution

Cannot model long-term dependency (due to the limited kernel size).

Naive Approaches v2

How about modeling just the next word?

The formal president of the United States is Donald <u>Trump</u>. Wherever Donald <u>Trump</u> went, security was tight.

The comic Donald <u>Duck</u> is very famous.

Context matters to the prediction!

"The comic Donald <u>Duck</u> is very famous."

Context Propagation

- The meaning of the sentence depends on the **context**.
- The context at $x_i^{< t-1>}$ is propagated to the next time step via the **hidden state** $h_i^{< t-1>}$.
- The next prediction is then modeled with both $x_i^{<t>}$ and $h_i^{< t-1>}$.
- The introduction of the hidden state h_i allows the network to have **memory**.

Recurrent Neural Networks

In recurrent neural networks (RNNs), weights are reused across multiple time steps.

RNN
$$\begin{aligned} h_t &= \sigma(Ux_t + Vh_{t-1} + b_h) \\ O_t &= \sigma(Wh_t + b_o) \end{aligned}$$

Non-RNN $O_t = \sigma(Wx_t + b)$

Back propagation through time

Back propagation through layers

Anatomy of an RNN Cell

 $h^{<t>} = \tanh(Ux^{<t>} + Vh^{<t-1>} + b_h)$

 $\hat{y}^{<t>} = \operatorname{softmax}(Wh^{<t>} + b_y)$

Hyperbolic Tangent as an Activation Function

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

Between **[-1, 1]**, the function behaves **linearly** (more or less).

As long as the input range is within [-1,1], the derivative in this range is approximately 1.

The function always scales the input to the wellcontrolled range of [-1,1].

Long-term Dependency

Sometimes, the contexts are far away. The <u>clouds</u> are in the <u>sky</u>. I grew up in the Netherlands... I speak fluent Dutch.

Contexts can be forgotten; the gradients are vanishing during back propagation.

Vanishing Gradient Problem in RNNs

Residual connection is a general solution to the vanishing gradient problem.

We should instead regulate which information to keep and which to forget.

However, for very long dependency (which is common in RNNs), this solution is not sufficient.

Long Short-term Memory

Hochreiter & Schmidhuber (1997)

LSTM regulates the flow of information through gates.

- **Input gate**: controls the extent to which a new value flows into the cell;
- Forget gate: controls the extent a value remains in the cell;
- Output gate: controls the extent to which the value in the cell is used to calculate the output activation.

$$\begin{split} f_t &= \sigma_g(W_f x_t + U_f h_{t-1} + b_f) \\ i_t &= \sigma_g(W_i x_t + U_i h_{t-1} + b_i) \\ o_t &= \sigma_g(W_o x_t + U_o h_{t-1} + b_o) \\ \tilde{c}_t &= \sigma_h(W_c x_t + U_c h_{t-1} + b_c) \\ c_t &= f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \\ h_t &= o_t \circ \sigma_h(c_t) \end{split}$$

Image source: Wikipedia

Gated Recurrent Unit

Cho et al. 2014 (arXiv:<u>1406.1078</u>)

GRU regulates the flow of information through gates:

- **Reset gate:** controls how much past information to forget;
- Update gate: controls what \bullet information to discard and what to keep (similar to the LSTM's input gate and forget gate combined)

$$z_{t} = \sigma_{g}(W_{z}x_{t} + U_{z}h_{t-1} + b_{z})$$

$$r_{t} = \sigma_{g}(W_{r}x_{t} + U_{r}h_{t-1} + b_{r})$$

$$h_{t} = z_{t} \odot h_{t-1} + (1 - z_{t}) \odot \phi_{h}(W_{h}x_{t} + U_{h}(r_{t} \odot h_{t-1}) - b_{t})$$

Generating Input Sequences

To apply a long sequence to a RNN, we need to define the training data x and the labels y:

We could either **reshape** it:

$$\begin{array}{c} \chi \\ (1, ...) (2, ...) (3, ...) (4, ...) (5, ...) (6, ...) (7, ...) \\ \chi \\ (8, ...) (9, ...) (10, ...) (11, ...) (12, ...) (13, ...) (14 \end{array}$$

Or truncate it using a **rolling window**:

$$(1, ...) (2, ...) (3, ...) (4, ...) (5, ...) (6, ...) (7, ...) (7, ...) (8$$

$$(3, ...) (4, ...) (5, ...) (6, ...) (7, ...) (8$$

$$(5, ...) (6, ...) (7, ...) (8$$

Stateful RNNs

A stateful RNN passes the **last state** of the **previous batch** as the **initial state** of the **next batch**.

This is useful if there is some connections between batches. This feature essentially allows long sequences to be **folded** into batches.

A non-stateful RNN initialises state of each batch to zero.

Batch 1	(1,) (2,) (3,) (4,) (5,) (6, .
Batch 2	(8,) (9,) (10,) (11,) (12,)

LSTM autoencoder for rare event detection

Input sequences

LSTM layers

Reconstructed Input sequences

Topology of Sequence Processing Models

RNN/Generation

RNN/Captioning

Image: Andrej Karpathy

Sequence-to-sequence Translation

Seq2seq model: Sutskever et al., 2014, Cho et al., 2014, Vaswani et al. 2017)

With attention mechanism, the decoder RNN selectively uses relevant data encoded in the context **vector** to generate the sequence.

Attention Mechanism

Bahdanau et al., 2015

Take Home Messages

Sequence modelling

... is a task that aims to solve the **dependency** between **tokens** in a sequence; ... requires passing information along the flow of the sequence; ... can be categorised into classification, regression, generation, and translation.

Recurrent neural networks (RNNs)

- ... backward propagate through **time**;

Miscellaneous

- the vanishing gradient problem.
- LSTM and GRU use gates to regulate the flow of information.

... are NN architectures that allow weights and states to be shared across multiple time steps;

... are usually slow to train because the time step dependency cannot be parallelised.

Long short-term memory (LSTM) and gated recurrent units (GRU) are RNNs that designed to mitigate

• The encoder-decoder architecture is particular useful in RNNs for sequence to sequence translation.

