

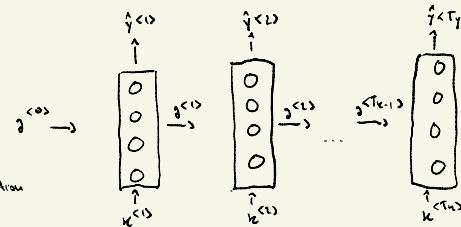

[RECURRENT NN]

Notation

RNN →

• Limitation: uses only info earlier in the sequence to make predictions.

BRNN solves this limitation



$$T_h = \text{size of } h \rightarrow T_h = T_y \text{ (only for this case)}$$

Forward propagation:

$$g^{<t>} = \vec{0}; W_{22}, W_{23} = \text{parameters Matrices}$$

$$g^{<t>} = g_L(W_{22}g^{<t-1>} + W_{23}g^{<t>} + b_2) \leftarrow \tanh / \text{relu}$$

$$\hat{y}^{<t>} = g_L(W_{23}g^{<t>} + b_3) \leftarrow \text{i.e. sigmoid}$$

$$g^{<t>} = g(W_{22}g^{<t-1>} + W_{23}g^{<t>} + b_2)$$

$$\hat{y}^{<t>} = g(W_{23}g^{<t>} + b_3)$$

$$g^{<t>} = g(W_{22}[g^{<t-1>} x^{<t>}] + b_2)$$

$$\hat{y}^{<t>} = g(W_{23}g^{<t>} + b_3)$$

$$\begin{bmatrix} W_{22} & W_{23} \\ [g^{<t-1>} x^{<t>}] & \end{bmatrix} = \begin{bmatrix} g^{<t-1>} \\ g^{<t>} \end{bmatrix}$$

Different types of RNN

$T_h = T_y$: many-to-many arch.

$T_h \neq T_y$: • Sentiment classification: $\begin{cases} h = \text{context} \\ y = 0/1, 0 \dots s \end{cases}$
many-to-one arch.

• Sequence generation: one-to-many arch.

• Machine translation: many-to-many, where $T_h \neq T_y$

• Sampling sequence model: allows you to generate new sequences from your RNN language model.

Sequence generation

• Language modelling:

→ Speech recognition: $P(\text{sentence}) = ?$
probability

$P(y^{<1>} | y^{<2>} | \dots | y^{<T_y>})$

Probability of sequence of words

→ Build a language model:

1) Training set: large corpus of English text

2) Test set: one-hot encoding, <EOS>, <UNK>

3) $y^{<1>} \rightarrow$ softmax with N outputs where N corresponds to
→ given the first 2 words are **code coverage**
What is the probability of **year**?
 $P(\cdot | y^{<1>} y^{<2>})$

→ A sequence model, models the choice of any particular sequence of words.

Issues with RNN

• Vanishing gradient (Vg), RNN fail to capture long term dependencies.

• " " " much harder to solve than exploding gradients.

Gated Recurrent Unit (GRU)

• Helps with Vg . it introduces a new cell, called C (memory cell) • $C^{<t>} = g^{<t>}$
it memorises old values in the sequence

LSTM

- Some purpose of GRU: allow you to learn very long sequences.

- Ruining equations are:

$$\hat{c}^{<t>} = \tanh(W_c [z^{<t-1>}, h^{<t>}] + b_c)$$

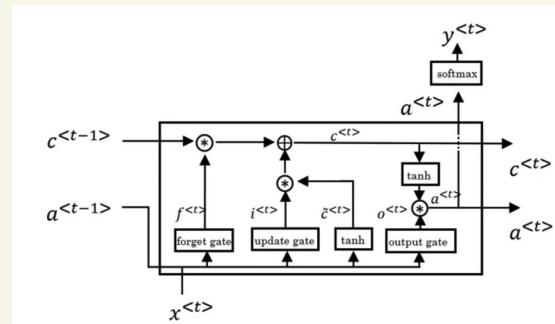
$$(\text{update}) \quad f_u = \sigma(W_u [z^{<t-1>}, h^{<t>}] + b_u)$$

$$(\text{forget gate}) \quad f_f = \sigma(W_f [z^{<t-1>}, h^{<t>}] + b_f)$$

$$(\text{output}) \quad f_o = \sigma(W_o [z^{<t-1>}, h^{<t>}] + b_o)$$

$$c^{<t>} = f_u * \hat{c}^{<t>} + f_f * c^{<t-1>}$$

$$h^{<t>} = f_o * \tanh(c^{<t>})$$



- Easier to build deeper NN with GRU. But LSTM is more used

Bidirectional RNN and Deep RNN

- BiRNN: getting information from the future

- Adding a backward layer for each forward layer