# 5- recurrent neural networks
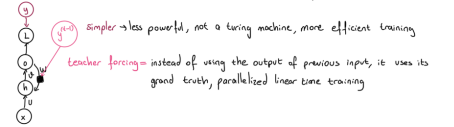
✴ typically for processing sequential data
✴ share parameters through time → otherwise cannot generalize to unknown sequence length
✴ have cycles/feedback connections

$$h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta)$$

hidden state (fixed length) — model — input (arbitrary length) — shared parameters

h ← means $h(t) \leftarrow h(t-1)$ one step back is forward

ex: I love apples → $h^{(3)} = f(f(f(h^{(0)}, x^{(1)}; \theta), x^{(2)}; \theta), x^{(3)}; \theta) = g^{(t)}(x^{(t)}, x^{(t-1)}, \dots x^{(1)}; \theta)$

✴ if we used this, diff param for each, and could work only with fixed len

(y) label sequence
(L) loss function
(o) output sequence
(h) — input sequence $V^t \to h^{(t)} \to o^{(t)}$
(x) input sequence

→ turing machine
$a^{(t)} = b + W h^{(t-1)} + U x^{(t)}$
$h^{(t)} = \tanh(a^{(t)})$
$o^{(t)} = c + V. h^{(t)}$
$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$

✴ the gradient of $L$ wrt $h^{(t)}$ depends on all the h's up to it
BPTT = back propagation through time
✴ cannot parallelize
✴ time and memory complexity = $O(t)$ len

✴ $U, W, V, b, c$ are learnable parameters
$L(\{x_1, x_2 \dots 3, \{y_1, y_2 \dots 3\}) = \sum_t L^{(t)}$
✴ there is an output for each input in the sequence

(y) — $\hat{y}^{(t-1)}$ Simpler → less powerful, not a turing machine, more efficient training
(L)
(o)
(h) teacher forcing = instead of using the output of previous input, it uses its ground truth, parallelized linear time training
(x)

bidirectional RNNs = output at t depends on both the past and the future
encoder-decoder sequence-to-sequence architecture =
↳ encoder RNN reads/processes the input sequence, then emits the learned content C
↳ decoder produces the output sequence based on C
vanishing or exploding gradient = weights will be $w^t$ at RNN with t layers/sequence
✴ either will explode (w>1) else vanish (w<1)
✴ gradient clipping = $\hat{g} \leftarrow \text{threshold}/\|g\| * \hat{g}$ for exploding gradient
✴ regularizing the gradient = for vanishing gradient

### long short-term memory (LSTM)



$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$
$\tilde{C}_t = \tanh(w_c[h_{t-1}, x_t] + b_c)$

$f_t$: forget gate (the percentage to remember)
$f_t = \sigma(w_f[h_{t-1}, x_t] + b_f)$
$i_t$: input gate (update the long term memory)
$i_t = \sigma(w_i[h_{t-1}, x_t] + b_i)$
$o_t$: output gate
$o_t = \sigma(w_o[h_{t-1}, x_t] + b_o)$

gated recurrent unit (GRU) = simpler LSTM
↳ the forget and input gate combined into a single update gate
↳ merges the cell and the hidden state (no $C_t$)

RNN    $O^{(t)}$            LSTM    $O^{(t)}$ → there is an output for every input



params $\theta = \{U, V, W\}$
state h
$h^{(t-1)}$

$\dots x^{(t-1)}, x^{(t)}, x^{(t+1)} \dots$

$h^{(t)} = \sigma(W h^{(t-1)} + U. x^{(t)})$
$o^{(t)} = \sigma(V h^{(t)})$

params $\theta = \{U, V, W, w_f, u_f, w_i, u_i\}$
state h
$h^{(t-1)}$

$\dots x^{(t-1)}, x^{(t)}, x^{(t+1)} \dots$

$\tilde{h}^{(t)} = \sigma(W. h^{(t-1)} + U x^{(t)})$
$h^{(t)} = f. h^{(t-1)} + i. \tilde{h}^{(t)}$
forget gate        input gate
$f = \sigma(w_f h^{(t-1)} + u_f x^{(t)})$    $i = \sigma(w_i h^{(t-1)} + u_i x^{(t)})$

RNN applications =
✴ character level language modeling = predicting next char in a given sequence
↳ one-hot encoding of characters
✴ word level language modeling =
↳ word embeddings = to represent each word with fixed-dimensions vectors
word2vec = widely used word-embedding model (2013)
↳ tends to map semantically similar words to nearby points
• (<content>, <word>) pairs: $x_j^t \approx t \dots \longrightarrow ([x_i, t], y), ([y, t], x)$. $N = \perp$ (btw)
• (biggest − big) + small = smallest (detects semantic relationships)
✴ image captioning = generating textual description of an image (image is an input too)
✴ machine translation = two RNNs (encoder & decoder), different sized input and output
↳ used for language translation and speech recognition
↳ encoder → summary (s) → decoder (performance decreases with longer sentences)
↳ attention mechanism: google's neural machine translation (transformers)
✴ pixel RNN = predicts the next pixels in an image, can also generate images from scratch
✴ neural turing machines (NTM) = increased memory size (RNN limited due to vanishing gradient)
↳ selective read or write on the memory
↳ every component is differentiable
↳ better than LSTMs in: copy, repeated copy, associated recall, sorting