

applications:

classifying sentiment
recognizing phonemes
classifying video activity

dealing with variable-size inputs \rightarrow zero-pad (simple)
 \rightarrow one input per layer
 $z^L = W^L a^{L-1} + b^L$ ($a^{L-1} = [x_{i,t}^L]$) \hookrightarrow each layer: activations of prev. layer inputs to layer.
have a "pretend layer" w/ $a^0 = 0$ to start chain.
RNN: W^L is the same for all layers, (and b^L).
 \hookrightarrow parameter-sharing, "extends a standard NN along the time dimension".
Minor backprop change: Instead of "overwriting" gradients at l for $l+1$, analogous to function chain rule. \rightarrow ADD to loss. $dZ/d\theta \leftarrow dZ/d\theta + d^l/d\theta$.
AKA "accumulate" the gradient.

applications:

generating text
Caption for image
predicting sequence of future video frames
generating audio sequence

dealing with variable-size outputs \rightarrow one output per layer
 \rightarrow separate loss for every layer
 $z^L = W^L a^{L-1} + b^L$
 $a^L = \sigma(z^L)$
 $\rightarrow \hat{y}_t = f(z^L_t)$
e.g. simply linear softmax
 $L(\hat{y}_1:T) = \sum_l L_l(\hat{y}_l)$
Reminder: $\delta_{x_f} = \frac{dL}{dx_f} = \frac{dL}{dy_f} \delta_{y_f}$ ($\delta_{y_f} = \frac{dL}{dy_f}$)
For variable-size output: $\delta_{y_f} = \delta_{y_f}^1 + \delta_{y_f}^2$ (for nodes with >1 dependents)
(separate losses, all but final layer activations)

Input and output at each step? combine the two concepts.

$\hat{y}_{t-1}, \hat{y}_t, \hat{y}_{t+1}, \hat{y}_{t+2}$
 $x_{t-1}, x_t, x_{t+1}, x_{t+2}$
Issues: RNNs are super deep, essentially, unable to retain memory.
For lots of layers, can run into Vanishing Gradient Problem
Could also have Explosive Gradient, could fix w/ clipping

"RNN dynamics"

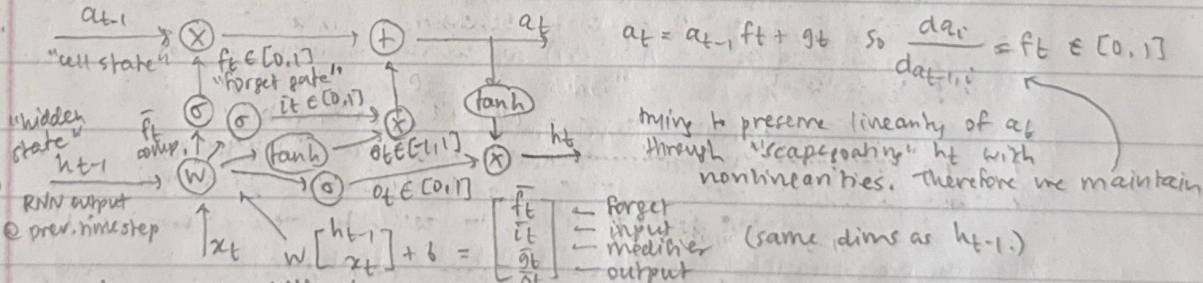
$a_t = \sigma(z_t)$ where $z_t = W a_{t-1} + b$
 $a_{t-1} = [a_{t-1}^1, \dots, a_{t-1}^n]$
 $a_t = \sigma(z_t)$ w/ Jacobian: $\frac{da_t}{da_{t-1}}$, should $\approx I$. (to remember)

to "remember" \rightarrow copy previous activation

to "forget" \rightarrow overwrite it w/ current input.

to and behold, LSTM.

(long short-term memory)



a_t : changes very little, long-term.

h_t : changes all the time, short-term.

most RNNs in practice look like this, because there are strong dependencies (w/ outputs like in text generation). Yesterday's output \rightarrow today's input!

Training ideas \rightarrow outputs are shifted inputs (as "tokens"). Better mistakes can compound...
Can use scheduled sampling to regularize

Model can see its own inputs and can act more reasonably.

Different ways to "use" RNNs?

one \rightarrow one \rightarrow stackable \rightarrow can incorporate continuous \rightarrow can be bidirectional (e.g. speech recognition).

RNN encoders \rightarrow process inputs prior to feeding them RNN layers (ex. images)

RNN decoders \rightarrow process outputs right after feeding them RNN layers (less common)

Cool examples:

\rightarrow Shakespeare text gen

\rightarrow Fake LaTeX gen for math

\rightarrow OpenAI GPT-2 gen text (using Transformer).