

114 - RNNs: stability Analysis LSTM

review

- Heated structures - short term dependencies
- Recurrent - " long term - "
- (*) This is not entirely true

(x) recurrent structures vs static structures

- (*) stage now symbolic ops included in MLPs
- (*) addition problem - 2 N bit numbers
- addition problem - 2 N bit numbers
- 'learn' this operation
- network will be large (exponential)
- training data will be exponentially large (and required)

(x) MLPs vs RNNs

- cast as recurrent problem advantages
- RNN: could compute using few neurons what a fixed/static structure would take an exponential no. to compute
- (*) (*) It's not merely problems with long term dependencies that are benefitted by RNN style architecture
- MLP with RNN architect can be more efficiently in terms of network form, amount of training data

(x) types of recursion

(x) Behavior of recursion

- finite response / time delay
 - If input is always bounded; what is condition for output to blow up?
 - (*) If output is a multiple of input that is greater than 1
(so input, $x_{t+1} = kx_t$ $k > 1$)
 - (*) CR: will never 'explode' (?)
- unless activations have singularities, output will remain bounded if input is bounded
- { (*) review; make

(*) Time delay structure
- bounded input, bounded output stability BIBO stability (P2)
SP term

(*) Is this BIBO?

①: what are conditions under which output will not explode?

If hidden activations and outputs are bounded

②: will it saturate and where? - saturate \rightarrow meaningless results.

what if activations linear?

BIBO stability: output carries info about input in a way that the isn't saturation

(*) Analysing recurrence

- output depends only on curr hidden activation

- If you want to know how output will behave \rightarrow analyse hidden layer h_R

(*) Examine linear activations

(*) Streetlight effect

- analyse linear systems; extrapolate to non-linear systems

$$z_R = w_h h_{R-1} + w_x x_R \quad h_R = z_R$$

(*) Non-linear systems

- unfold the recurrence

(B): Review - expand recursion

$$h_R = w_h^{R+1} h_{-1} + w_h^R w_x x_0 + w_h^{R-1} w_x x_1 + w_h^{R-2} w_x x_2 + \dots$$

$$h_R = H_R(h_{-1}) + H_R(x_0) + H_R(x_1) + H_R(x_2) + \dots$$

(i): - response to an input x_0 at time 0; when there are no other inputs, zero initial condition

(e): rewrite operations: $H_R(x_0)$ - response at time R assuming input x_0 at time 0, assuming all other inputs 0. (ceteris paribus)

(*) As system is linear; can consider x_0 as scaling factor; pull it out yielding:-

$$h_R = h_{-1} H_R(l_{-1}) + x_0 H_R(l_0) + \dots + x_L H_R(l_L) + \dots$$

(*) for linear activation; responses at time k can be thought of as summations of individual inputs, assuming all else 0.

(*) examine what happens given network receives single input at time 0.

- If that blows up \rightarrow system blows up
- If that does not blow up \rightarrow \dots not blow up
- principle of superposition in linear system A4 Review

(1) scalar recursion

- scalar recursive

$$h(t) = w h(t-1) + c x(t)$$

$$h_0(t) = w^t x(0)$$

- subscript 0 means
input occurred at time 0 !

- behaviour of $h(t)$?

Q: Similar to argument previously; $|w| > 1 \quad t > 1 \quad \text{if}$

- stability is only when $w=1$.

(2) vector recursion: vector recursion

$$h(t) = W h(t-1) + C x(t)$$

$$h_0(t) = W^t x(0)$$

- use eigen decomposition of W (weights matrix). $W = U \Lambda U^{-1}$

BR: columns of U ? (eigenvectors)

Q: - orthonormal (basis vectors)

- linearly independent

- note $W_{ui} = \Lambda_{ui}$

for any vector $x' = cx$

← (AS) slides S.2020
→ part

$$x' = a_1 u_1 + a_2 u_2 + \dots + a_n u_n$$

$$Wx' = w_{11}u_1 + w_{21}u_2 + \dots + w_{n1}u_n$$

$$= a_1 w_{11} + a_2 w_{21} + \dots + a_n w_{n1} \quad (\text{as } a_i \text{ is scalar})$$

$$\overset{\sim}{= \lambda u_1} \quad \overset{\sim}{= \lambda u_2}$$

$$Wx' = a_1 \lambda_1 u_1 + a_2 \lambda_2 u_2 + \dots + a_n \lambda_n u_n$$

$$W^t x' = a_1 \lambda_1^t u_1 + \dots + a_n \lambda_n^t u_n$$

(*) Consider behavior as $t \rightarrow \infty$

$$\lim_{t \rightarrow \infty} |W^t x'| = a_m \lambda_m^t u_m = a_m \max_j |\lambda_j|$$

(a6): BR's arguments are a little unclear to me here
(convergence)

(@): ability of network to 'remember' depends on largest eigenvalue of W.

(*) 'all of other a_i 's' disappeared

- (?) notice what is meant by length of hidden vector will expand/contract

(@): $|\lambda_{\max}| > 1$? - determines blow up

(*) complex eigenvalues \rightarrow oscillation + blow up.

Lesson...

(@): linear systems:-

long term behavior depends on eigenvalues of the hidden weights matrix. $|\lambda_{\max}| > 1$

- largest eigenvalue $> 1 \rightarrow$ system blows up

- $-\frac{1}{|\lambda_{\max}|} < 1 \rightarrow$ response will vanish quickly

- complex eigenvalues cause oscillatory behaviour on top of divergence/blow-up

- smooth behaviour \rightarrow force weights matrix to have real eigenvalues
enforce

- symmetric weights matrix

(A6)

(*) How about non-linearities (scalar)

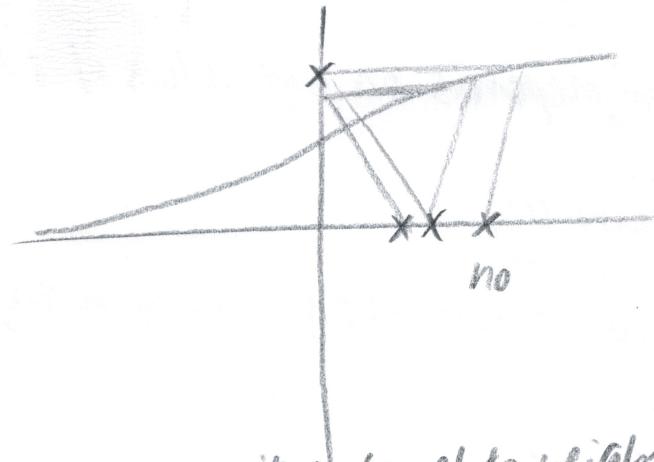
- consider :-

$$n(t) = f(w_h(t-1) + c x(t))$$

$$h_t = \sigma(w_h t_1 + x)$$

Q: what is largest slope of sigmoid function?

- "worst"
- assume it's 1
- what happens when initial value at h_0



(**) Some key insights on behavior of non-linearities (in rel to weights, biases)

- sigmoid, tanh, ReLU

- diagrams on slide ① and ② differ by whether initial value \hat{h}_0 is to the left of 0 (< 0) or > 0 on the plot of output vs input on diagram above.

(*) vector process

- analysis more complex

$$n(t) = f(w_h(t-1) + c x(t))$$

- uniform unit vector init: $[1, 1, \dots] / \sqrt{N}$

(*) Behavior similar to scalar recursion

(*) Stability Analysis

- need more specialised tools for considerations of stability formally
i.e. Lyapunov stability, Routh's criterion

(*) (**) Comments on saturation profiles:- (activations all eventually saturate)

i) sigmoid: saturate quickly

ii) tanh: retain memory about where you started from (initial cond?)
for a long period of time (but can eventually saturate)

- (*) ReLU → not used (blow up)
- ② (67). Activations in RNNs tend to be stuck - retain 'memory' for long enough
- (*) Behaviour ^{also} depends on eigenvalues of W .
-
- (*) Deep recursion
- OR: activations do not change significantly
-
- (*) Story so far ②
- excellent summary
-
- (*) Vanishing gradients
- (*) Issue with any deep network
- gradient of error/loss at beginning of network is 'unstable'
 - why?
- (67): note nested function rep. of NN
- MUL: $y = f_N(W_{N-1}f_{N-1}(W_{N-2}f_{N-2}(\dots(W_0x)\dots))$
- W_k - weights mat
in k^{th} layer
- Error/Div: $\text{Div}(x) = \partial(f_N(W_{N-1}f_{N-1}(W_{N-2}f_{N-2}(\dots(W_0x)\dots)))$
-
- (*) Training deep nets
- How to find $\nabla_x f$ where $f = f(Wg(x))$
 - $z = Wg(x)$
 - $\nabla_x f = \nabla_z f \cdot \nabla_g z \cdot \nabla_x g = \nabla_z f \cdot W \cdot \nabla_x g$
- $\nabla_z f$ - Jacobian matrix
of $f(z)$ wrt z
i.e. $J_z(f)$

for

$$\text{div}(x) = D(f_N(w_{N-1}f_{N-1}(w_{N-2}f_{N-2}(\dots(w_0x))\dots)) \quad (\text{should be } \text{div}(x, y))$$

- y suppressed

$$\nabla_{f_k} \text{div} = \nabla D \cdot \nabla f_N \cdot w_{N-1} \cdot \nabla f_{N-1} \cdot w_{N-2} \dots \nabla f_{k+1} \cdot w_k$$

- $\nabla_{f_k} \text{div}$ - gradient of error $\text{div}(x)$ wrt output of k^{th} layer of network
 - required to compute $\nabla_{w_{k-1}} \text{div}$ and $\nabla_{b_{k-1}} \text{div}$.

- ∇_{f_k} is Jacobian of $f_k()$ wrt current input

- clarity:-

(*) Jacobian of hidden layers (RNN)



$$h_i^{(l)}(t) = f_l(z_i^{(l)}(t))$$

$$\nabla f_l(z_i) = \begin{bmatrix} f'_{l,1}(z_1) & 0 & \dots & 0 \\ \vdots & f'_{l,2}(z_2) & & \vdots \\ 0 & 0 & \dots & f'_{l,N}(z_N) \end{bmatrix}$$

(*) derivative (subgradient) of activation is always bounded.

- diagonals/singular values of Jacob - bounded.

(*) limit on how much multiplying a vector by Jacobian will scale it.

Notation: $\nabla f_l()$ is the derivative of output of the (layer of) hidden recurrent neurons with respect to their input

S. 2019:-

- vector activations - full matrix

- scalar activations - matrix where diagonal entries are derivatives of hidden layer activations

(**) - more about masking

- derivative of hidden state act.

- most common activation fns : sigmoid(), tanh, ReLU

have derivatives less than 1.

(*) most common RNN act (ReLU) - tanh()

- derivative of tanh() is never greater than 1!

(A8) - check
on Wolfram

(6D): Multiplication by Jacobian \rightarrow shrinking operation

(7D): Backprop \rightarrow Jacobians shrink the derivative.

(*) what about weights

(6D): single layer RNN - identical weight matrices (A9) (?) singular values?

can product $\nabla_{\text{fw}} \text{Div}$:

i) expands 1D along directions in which singular values of weight matrices are greater than 1.

ii) shrinks 1D in directions where singular values are less than 1. \downarrow SVD constraints.

- does he mean eigenvalues of W ? (rather than singular values?)

(*) Hence exploding/vanishing gradients

(*) gradient problems in deep nets.

- Hence gradients of div/error/loss wrt to earlier layers in network (i.e. further 'back' in backprop, from greater no. of multiplications) can explode/vanish.

- Has effect on gradient descent updates, ↑ problem with net. depth.

(*) vanishing gradient examples

(Q): At what stage should derivatives be largest?

(A): Increasing n means derivatives of loss wrt parameters

- If so; will be largest at beginning of training ✓

(*) Illustration of vanishing gradients

- true for generic neural networks

(*) Story so far

- key points - review

(*) Recurrent nets are deep nets

- even though vanishing/exploding gradients apply to deep nets (rather than RNNs specifically); they apply to RNNs if you consider these as an instance of deep nets

or: we want RNNs to have long memory (impaired by vanishing/exp. gradient problems).

(*) Can happen on forward pass too.

(*) Long-term dependency problem

- we want a way of 'remembering'; but we are in a situation where the extent of this (i.e. lack of saturation) depends on (eigenvalues of) network params. (rather than inputs)

↳ LSTM - mystery/magic?

(*) Exploding/vanishing gradients

- recall: exploding/vanishing gradients is antithetical to 'memory retention'

(*) Exploding/vanishing gradients:-

- i) Behavior of Jacobians J_{f^n} wrt current input

(Bands
on J_{f^n}) depends on nature of activation

- ii) Eigenvalues of weight matrices W_K

- not related to input

- (iii): whether network 'remembers' depends on parameters (w) rather than what it is 'trying to remember' - input.

- (*) looking for a different design (a failing of RNN).
 - input-based determination of 'memory'.
- (*) use example of curly brace in sentence being remembered
- (A7) design choice/aspiration:-
 seek a design/architectural net:-
- i) doesn't cause exploding/vanishing gradients
 - ii) 'retains' useful memory arbitrarily long (^{should now} what to remember based on input)
- $\text{memory}(k) \approx c(x) \cdot o_k(x) \cdot o_{k-1}(x) \dots o_1(x)$
 - $\nabla c(x) \nabla o_k(x) \nabla o_{k-1}(x) \dots \nabla o_1(x)$
 - info extracted from input
 - determination of whether to remember - based on input based flag (by a function of input)
 - no weights, no activations; only tiggers that bump up/down memory.

(*) constant error cascade

- history stored in variable $c(t)$
 - no weights, no non-linearities
- (*) non linearities done by other portions of net.
- $\sigma(\cdot)$ is the gate: depends on want input, hidden state. (and other stuff)
- curly brace example:- (e.g. forget open brace)
- i) seen closed brace - current state?
 - ii) know open brace has been detected and stored. - hidden state.

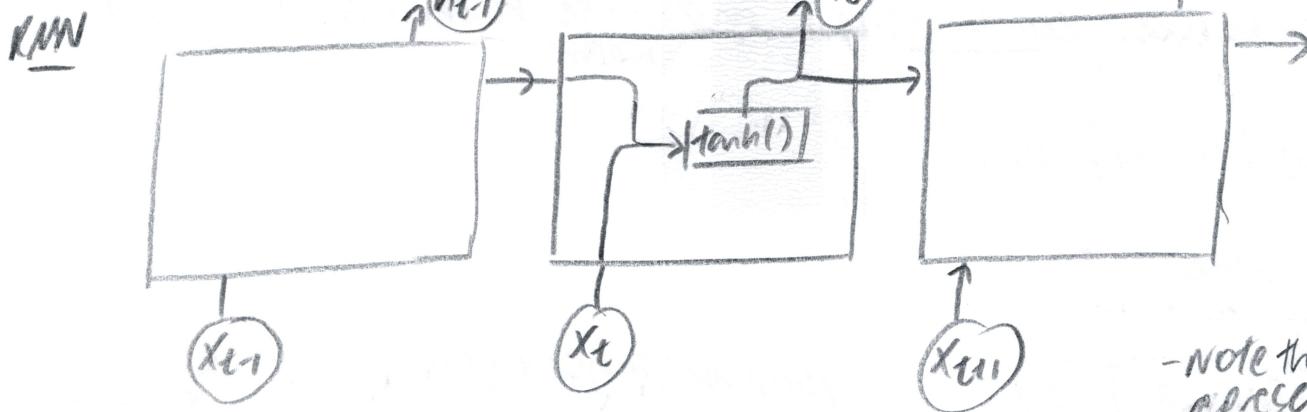
(A8) - review architecture

(A8) - issue with indices
 in 5.2020 slides of CEC diagram.

(*) enter LSTM

- (*) the key point is that it is an input-based detem.

(*) note representations of RNNs / LSTMs :-

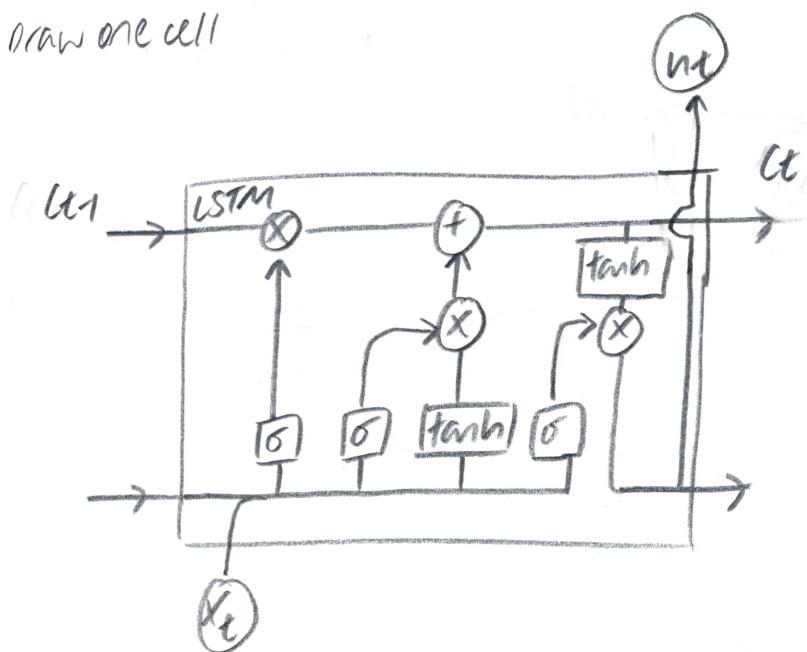


- note the representation established at beginning.

- LSTMs :- Quite complex 'circuit diagrams'

- draw / pick up in review **AIO**

LSTM :- draw one cell



BR: breaks each element / functional sub-circuit of this down

- LSTM CEC: C - linear history carried by CEC

- carries info through, only affected by gate
- And addition of history (which is gated)

- LSTM: gates - Gates are simple sigmoid with output range (0,1)
- controls how much info to be let through.

LSTM: forget gates

- determines whether to carry our history or forget it
(*) distinction between cell memory C and state h_t that is carrying over the hidden-

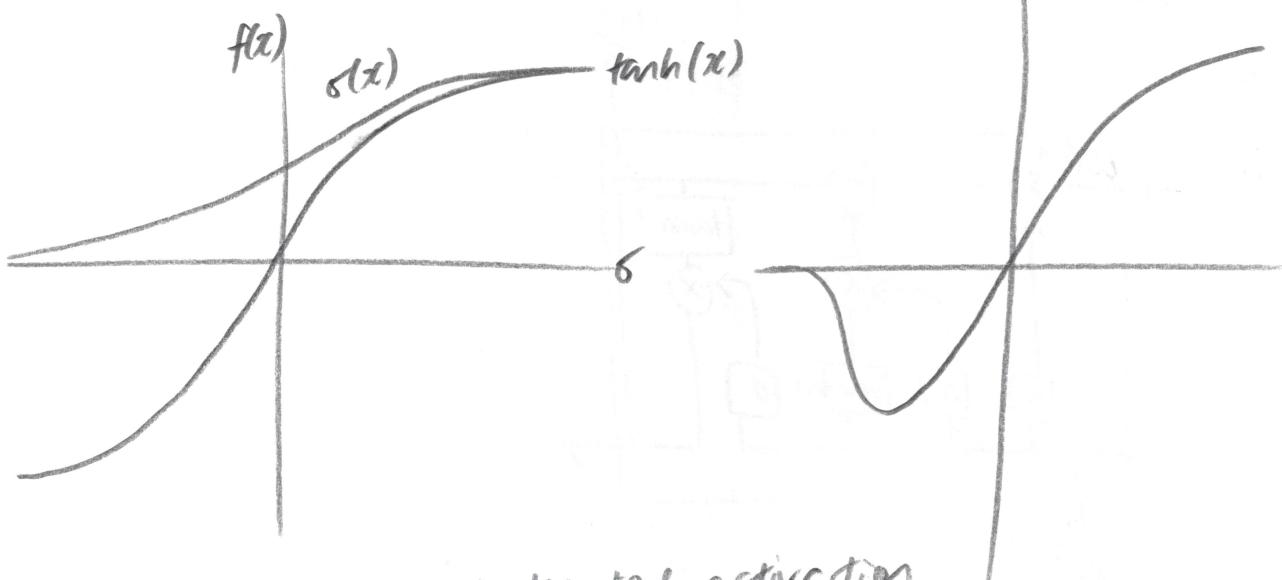
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

(*) LSTM: Input Gate

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
 - 'perception' layer that dec. if sth new/int in input.

$$\hat{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$
 - a gate that dec.

product of σ and \tanh (math. view) :-



(*) can also view as more complicated activation
(so if going for mathemat. view)

(*) Memory cell update

$$c_t = f_t * c_{t-1} + i_t * \hat{c}_t$$

(*) previous memory (c_{t-1}) multiplied by output of forget gate f_t ;
added to product of input gate and pattern detector in input.

?) which is i_t and \hat{c}_t \textcircled{AII}

LSTM: Output and output gate

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o)$$

$$n_t = o_t * \tanh(c_t)$$

- input memory c_t goes through, passed by $\tanh()$ to compute hidden value n_t
- hidden state h_t also multiplied by a gate that 'decides' if memory contents are worth remembering
- (*) make a distinction between intuition/expl./naive; and maths.

LSTM: Peephole connection

(see slides)

- raw memory information by itself; can be input

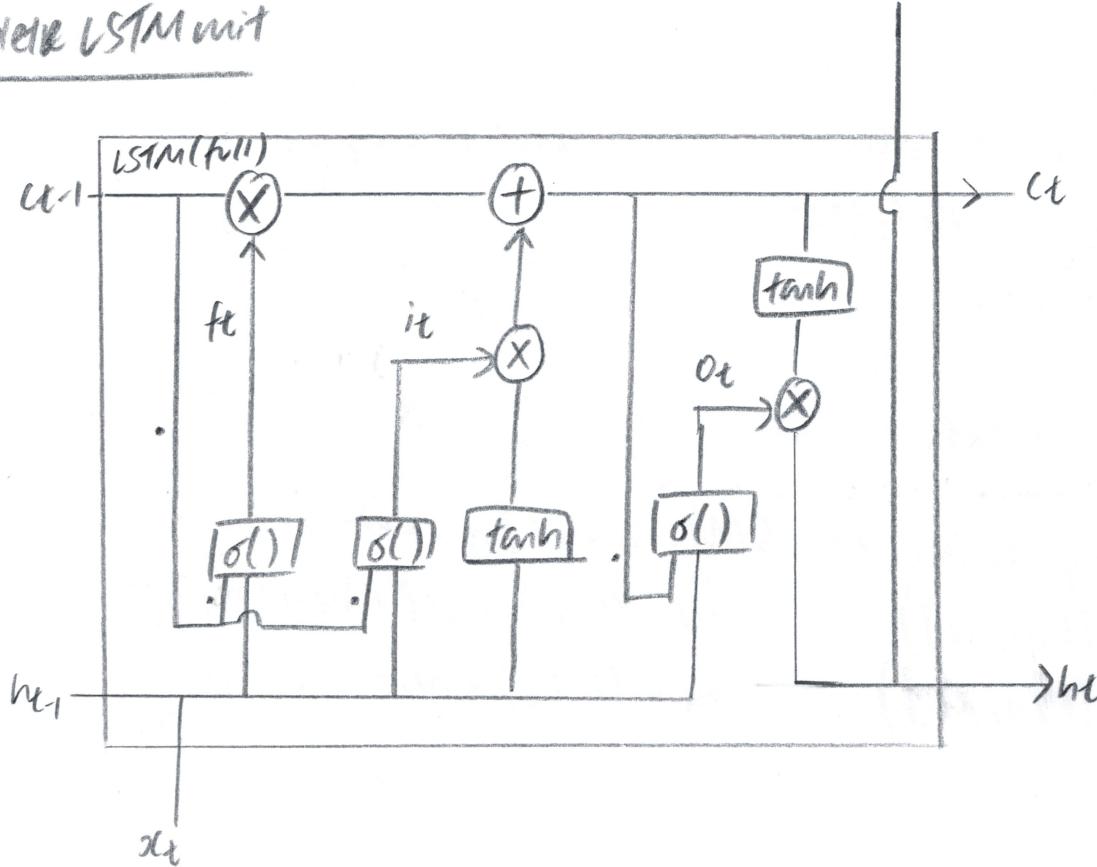
$$f_t = \sigma(W_f \cdot [c_{t-1}, h_{t-1}, x_t] + b_f)$$

(.) indicates
on diagram
peephole conn.

$$i_t = \sigma(W_i \cdot [c_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [c_t, h_{t-1}, x_t] + b_o)$$

(*) complete LSTM unit



Forward rules

$$f_t = \sigma(W_f \cdot [h_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [h_t, h_{t-1}, x_t] + b_o)$$

GATES

VARIABLES

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

$$c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

$$h_t = o_t * \tanh(c_t)$$

(*) LSTM forward pseudocode

- A12) - review pseudocode (trivial if you have equations mastered)

(*) Backprop rules: Backward

BR: Head can start spinning when doing this.

- A11) BR seems to recommend foregoing maths! And instead view directly in terms of code

- new to me.

(*) LSTM backward

- assume already have dh_0

- see how to compute
a deriv.

- A13) Review backprop with (LSTM)
suppl. slides

(*) GRUS: Simplify LSTMs

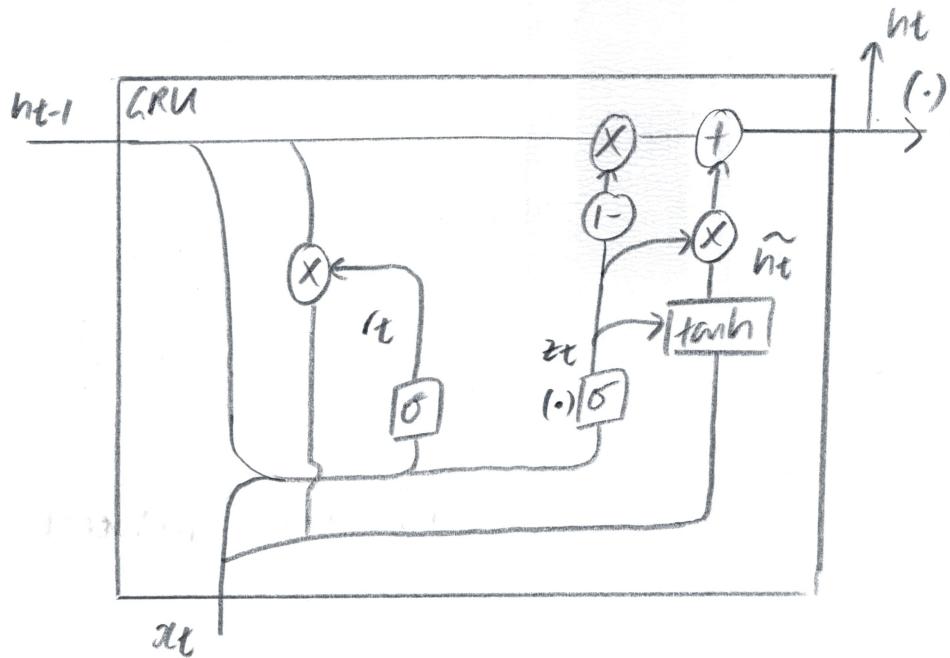
- redundant computation in LSTM: simplify via \rightarrow GRU

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$



(.) combine forget and input gates

If a new input is to be remembered, old memory forgotten

(.) don't separately store compressed, regular memories.

LSTM eq.

(14) - review

LSTM arch, Bidirect. LSTM

(15) - review.

