

L13-RNNs

(*) simple RNNs generate data at beginning

BR: Raises the idea of coding constraints; ie RNN generating output that works within constraints.

(*) Predictions are contingent on past predictions.

(*) Modelling series

- focus on sequential data

- stock example: vector of information / r.v.s. changing day by day.

(*) distinction between classification, prediction

(*) output may be series

(*) Representational shortcut

- use a diagrammatic representation of RNNs (different from before).

- important as we will use all throughout presentation

- A square represents series of components/neurons. (entire layer with many neurons)

- Arrows → full interconnections

- BR: Representation shows 2 hidden layers

(*) show what happens also when using 2 input vectors.

- Blue squares - inputs ('into page')

- Green squares - layer of neurons

- Arrows - complete set of interconnections

(*) Stock prediction

BR: examine all history of stock price to make current prediction.

- this is pedagogical; there are many assumptions/caveats with such an approach...

- CNN applied to series data → time delay NN.

(*) Stock predictor

- Sliding window

- There is an issue

(*) Finite response system

$$y_t = f(\underbrace{x_t, x_{t-1}, \dots, x_{t-N}}_{N\text{-width}})$$

- prediction at time t is function of previous N inputs.
- (A1): How can you contextualise this with respect to your existing knowledge of time series
- only 'finite window' of history considered.
- (B1) frames this in terms of history
- (B2) increase amount of 'history' being used to make predictions - increase window width. (N)
- BR: Does this actually satisfy what we require?

(*) Long-term dependencies

- note this is a signal processing instructor commenting on stock price movements; not econ/finance professor.
- interesting to see their view of the problem
- (W): A lot of assumptions being bundled into this example.

(*) Infinite response systems

- BR: we want to look at very long term dependencies "infinite memory"
- and do this without exponential blowing up of connections?

BR: what kind of feature of NN are we looking for?

(W): Memory \times Feedback (crucial cybernetic principle you didn't think about)

$$- y_t = f(x_t, y_{t-1})$$

- init state: $t=0 \rightarrow y_{-1}$

$$y_0 = f(x_0, y_{-1}); y_1 = f(x_1, y_0) \dots$$

- recurrence

- (*) a single output influences output 'for rest of time'
- (*) NARX network - nonlinear autoregressive network with exogenous input

- information about past (x_t, x_{t-1}, \dots) stored in output!

$$y_t = f(x_{0:t}, y_{0:t-1})$$

- Output contains info about entire past

(*) Satisfy ass. here; how do RNNs perform against time series and econometric models?

(*) We may also want past influence to dampen influence on present also...

(*) one-tap NARX network

- NARX net with recursion from output.
- entire sequence of computation shown (unfolded) (but still one NN)
- Brown boxes show output nodes $\exists \odot$
- Yellow boxes outputs
- similarity to scanning MLP (one NN)

(*) Difference between this and previous NNs covered?

- grey column in this representation is identical!

(*) Condense the yellow code \rightarrow brown nodes - output.

(*) Info is stored in output

(*) generic NARX network

(*) can use arbitrary units to feed into next.

$$y_t = f(x_t, \dots, x_{t-1}, y_{t-1}, \dots, y_{t-K})$$

- output at t , y_t depends on past K outputs and L inputs.

(*) complete NARX network (infinite past outputs, inputs) \rightarrow not practical...

(*) NARX network

- was popular for time-series (e.g. stocks, weather)
for phenomena with distinct seq. of innovations.

OK: Issue:- No genuine 'memory'; a 'substitute' is the 'storage' of info in y output (and not on the network).

(*) Nothing on network stores information (or dedicated to storage)
of info.

(*) Make memory explicit

(*) Introduce a memory variable

$$m_t = r(y_{t-1}, h_{t-1}, m_{t-1})$$

$$h_t = f(x_t, m_t)$$

$$y_t = g(h_t)$$

m - memory variable (stored in memory unit)

y_{t-1} - output ($t-1$)

h_{t-1} - hidden state ($t-1$)

m_{t-1} - " "

(*) Jordan Network

(*) memory unit - dark blue?

- introduced memory variable as described before

- introduce a running average of past output (memory unit)

(*) memory unit influenced hidden state at next time point.

(*) memory was 'fixed structure'; does not 'learn' to remember.

(*) Elman network

- separate memory state from output

- memory 'clones' previous state

- training \rightarrow slices network

- treat each slice as independent predictors with (?)

- use slices to estimate parameters

(*) no transfer of gradients from one to another (13) - review

(*) There are slides in S.2020

(NARX, Elman, Jordan networks) missing from S.2019

(*) Recurrence for Jordan, Elman \rightarrow partial

- only occurs during operational phase; not exploited during training

(*) Alternate model for finite response systems

- state space model (connection with 10-708?)

$$h_t = f(x_t, h_{t-1})$$

$$y_t = g(h_t)$$

ht - state of network

- model directly embeds
the memory in the state.

- need to define initial state h_{-1}

(*) Need to
define
an initial
hidden state
 h_{-1}

(*) h_t (state) - function of current input, past state h_{t-1}
(x_t)

y_t (output) - function of current state (h_t)

h_{-1}
as a
'parameter'

(*) distinction between this and previous? (see diagrams)

(*) State at any time is function of state at previous time.

↳ information 'passed forward'

(*) Fully recurrent, if modelling f and g as NNs
functions

(*) From now on; use recurrent as qualifier

(*) Simple state space model

(*) Curvilinear → cannot draw each column independently of every other column.

- information flowing through continuously.

(*) Have to show entire sequence of computations

(*) h has been used as a generic representation of a state

- no specification of what this is so far.

(*) 'State' can be arbitrarily complex

(*) single/multi-layer RNN / multiple recurrent layer RNN

(*) note architecture

- distinction - state defined in a different way.
- abstract: an skip connections

(*) A more complex state

- connections more complex; basic idea same
- (*) essentially showing variants on some principle

(*) or network may be even more comp!

- no longer a standard state space model
- AS there is recurrence in output - SSM - no recurrence in output

(*) generalisation to other recursions

- can generalise to recurrence over wider window (?)
- BR terms this 'deeper' recurrence; and I don't like the imprecision of the semantics
- via deeper recurrence \rightarrow initial states (further back) need to be specified e.g. 2 steps back h_1 h_2 for 2-step recurrence

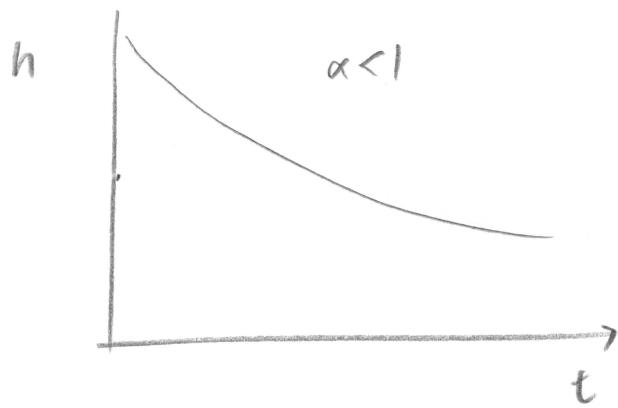
Illustrate

$$\text{E.g. } h_t = \alpha h_{t-1}$$

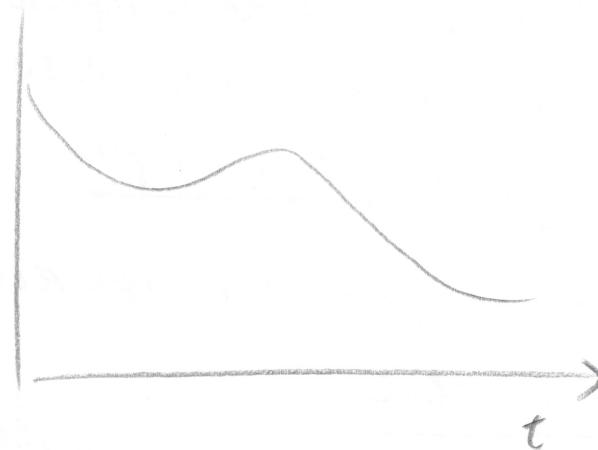
$\alpha < 1 \rightarrow$ decays exponentially
 h_t as $t \rightarrow \infty$

- suppose you want oscillation \rightarrow needs 'deeper steps'

$$h_t = \alpha h_{t-1}$$



$$h_t = \alpha h_{t-1} + \beta h_{t-2}$$



(*) "Deep recurrence"
can capture more complex patterns in memory

(*) simplest structures are most popular

(*) Recurrent Neural Network

(*) introduces 'loop notation' for drawing RNNs (similar to plate notation)
(*) feedback loop/ does not refer to a violation of causality
recurrence (there is temporal unfolding)

(*) feedback is to same variable in next period time.

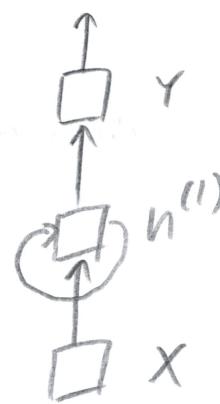
(*) Tailed/simplified

(*) Enfolded and unfolded representation of RNN

(*) note \rightarrow squares - full layer
arrows - full connects. (e.g. $K \times L$ weight matrix)

(*) RNN equations (*)

Q:-



$h_i^{(-1)}$ = part of net. params.

$$h_i^{(1)}(t) = f_1 \left(\sum_j w_{ji}^{(1)} x_j(t) + \sum_j w_{ji}^{(1)} h_i^{(-1)} + b_i^{(1)} \right)$$

$$y(t) = f_2 \left(\sum_j w_{jk}^{(2)} h_j^{(1)} + b_k^{(2)} \quad k=1, \dots, M \right)$$

- superscript \rightarrow layer of net.
from which inputs obtained

i) current weights

- vector function at output

ii) recurrent weights.

- state/node activ - $f_i(h)$ - typically $\tanh()$

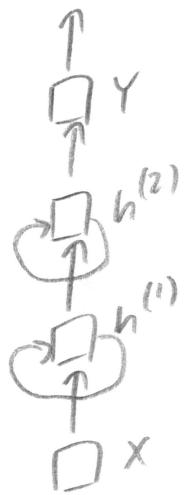
- every neuron - bias input

(*) Note activation $f_i()$ acts on a weighted sum of inputs AND

(*) $w^{(1)}$ - connection from neuron in 1st layer
'back' to neuron in 1st layer

previous time states
(with bias)

(*) Multiple recurrent RNN



$n_i^{(1)}(-1)$ = part of net params.

$h_i^{(2)}(-1) = \dots$

$$h_i^{(1)}(t) = f_1 \left(\sum_j w_{ji}^{(0)} x_j^{(t)} + \sum_j w_{ji}^{(1)} h_j^{(1)}(t-1) + b_i^{(1)} \right)$$

$$h_i^{(2)}(t) = f_2 \left(\sum_j w_{ji}^{(1)} h_j^{(1)}(t) + \sum_j w_{ji}^{(2)} h_j^{(2)}(t-1) + b_i^{(2)} \right)$$

$$Y(t) = f_3 \left(\sum_j w_{jk}^{(2)} h_j^{(2)}(t) + b_k^{(3)}, k=1, \dots, M \right)$$

- output - vector fn

- state node activ - fn() - typically tanh.

(*) Multiple RNN (skip conn.)



$n_i^{(1)}(-1)$ = network param.

$h_i^{(2)}(-1) = \dots$

$$h_i^{(1)}(t) = f_1 \left(\sum_j w_{ji}^{(0,1)} x_j^{(t)} + \sum_i w_{ii}^{(1,1)} h_i^{(1)}(t-1) + b_i^{(1)} \right)$$

$$h_i^{(2)}(t) = f_2 \left(\sum_j w_{ji}^{(1,2)} h_j^{(1)}(t) + \sum_j w_{ji}^{(0,2)} x_j^{(t)} + \sum_i w_{ii}^{(2,2)} h_i^{(2)}(t-1) + b_i^{(2)} \right)$$

$$Y(t) = f_3 \left(\sum_j w_{jk}^{(2)} h_j^{(2)}(t) + \sum_j w_{jk}^{(1,3)} h_j^{(1)}(t) + b_k^{(3)} \right) k=1, \dots, M$$

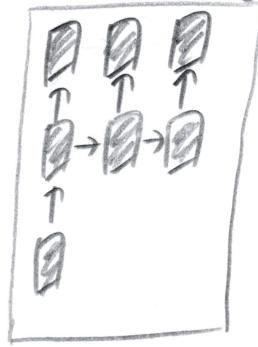
(*) Variants of RNN :-

1. conventional MLP / NN
2. sequence generation (single input, multiple outputs)
3. sequence prediction / class-

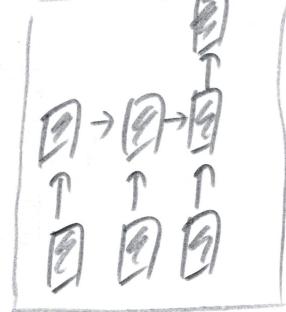
1. 1-to-1



2. 1-to-many



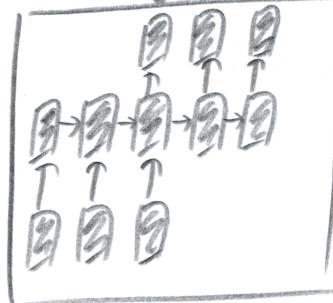
3. many-to-1



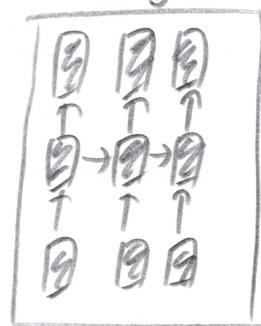
(*) variants.

1. delayed sequence
2. sequence to sequence

1. many-to-many



2. many-to-many



(*) Story so far

- (*) make sure you appreciate distinctions between NARX, 'simple' RNNs, SSMS.

(*) How to train network

- use BPTT - backprop through time.

- given a collection of sequence inputs

- (X_i, d_i) where $X_i = X_{i,0}, \dots, X_{i,T}$

$d_i = d_{i,0}, \dots, d_{i,T}$

(*) Train network params to minimise error between output of network

$y_i = y_{i,0}, \dots, y_{i,T}$ and desired outputs

(*) Assume synchronous output for every input.

(*) At the i th input consists of T input vectors (for each timestep) $x_{i,t}$
ith output ——— T output vectors — “ — $y_{i,t}$

(*) find a way of taking/estimating params to do this.

(*) Forward pass - RNN - asyn.

- fairly trivial

Ⓐ - review pseudocode (quickly).

(*) BPTT

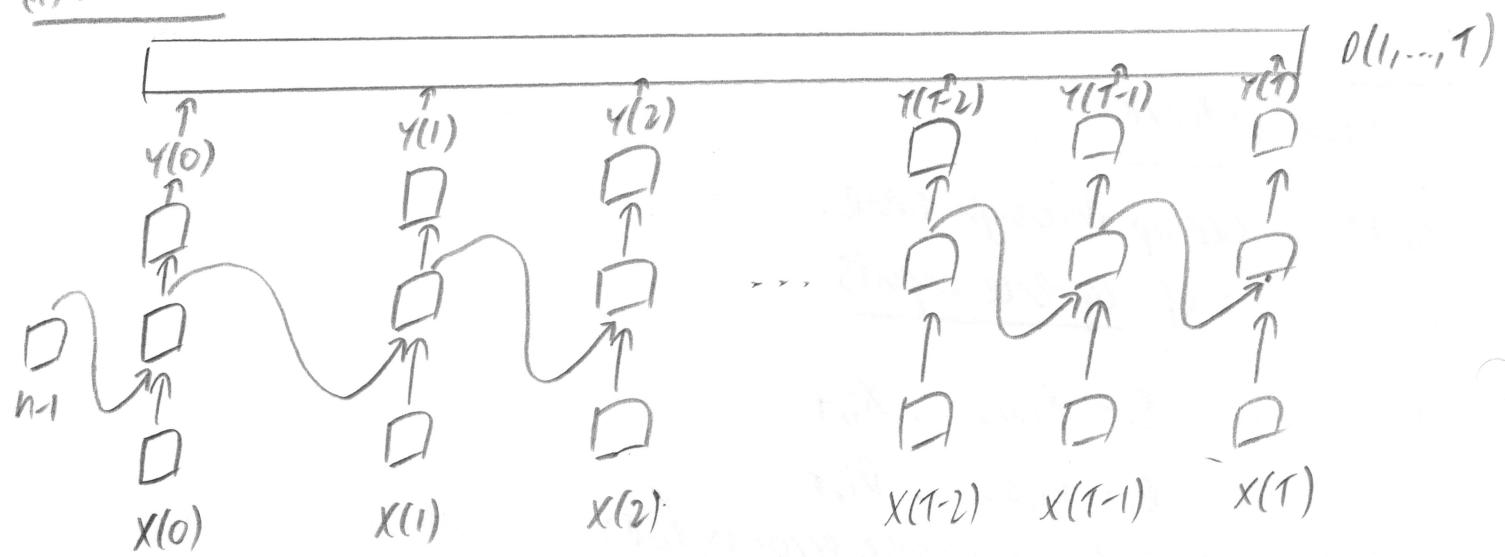
(*) Divergence is between entire sequence of outputs produced on forward pass and entire sequence of desired outputs.

Ⓐ and Ⓑ point is there is no reason to expect an additive relationship between individual divergence at each time point; and divergence over whole sequence.

(*) Due to alignments / subst. etc.

(*) But can be explicitly assumed in special cases.

(*) Notation



• $y(t)$ - output at time t

(one hidden layer example)

• $y_i(t)$ - i th output

• $z^{(2)}(t)$ - pre-activation value of neurons at output layer at time t

• $h(t)$ - output of hidden layer at time t

• $z^{(1)}(t)$ - pre-activation value of hidden layer at time t

Ⓐ - difference in subscript between S.2019/S.2020 slides

(*) BPTT

- (*) Note divergence is a vector function of output
- (*) Have to compute in totality
 - But do not confuse divergence being a vector fn with being able to take derivatives wrt a component of output (?)
 - ∂V is a function of $y(0), y(1), \dots, y(T)$
- (*) Compute $\frac{\partial D}{\partial y_i(t)} \forall i, t$ (i.e. overall $i=1, \dots, N$ observation sequences) use 'DIV' vs 'div' in previous to emphasise diff
 $t=1, \dots, T$ time steps
- (*) Can be a source of difficulty;

(*) BPTT

- (*) Explicitly assume that :-
overall divergence is an additive sum of local divergences at each time.
- (*) Require: $\frac{\partial D}{\partial y_i(t)} \forall i, t$
- (*) Will usually get $\frac{\partial D}{\partial y_i(t)} = \frac{\partial \text{Div}(t)}{\partial y_i(t)}$ (?)
- (*) Assumed for pedagogical purposes; in practice, BPTT is not contingent on this assumption.
- (?) Note picture showing this

(*) BPTT (equations)

compute $\frac{\partial \text{div}}{\partial y_i(T)}$ for all i

(for one training instance)

BR: uses this convention

(*) make sure these are crystal clear

- Assuming y is column vector; $\nabla_{y_t} \text{div}$ is a row vector

- conventional backprop

$$(*) \nabla_{Y(T)} \text{DIV} = [\quad]$$

$$\nabla_{Z^{(1)}(T)} \text{DIV} = [\quad] = [\quad]$$

⑧ $\nabla_{Z^{(2)}(T)} \text{DIV} = \nabla_{Y(T)} \text{DIV} \nabla_{Z^{(1)}(T)} Y(T) = \nabla_{Y(T)} \text{DIV} [\quad]_{Z^{(2)}(T)} (Y(T))$

(*) Note in scalar form:-

$$\frac{\partial \text{DIV}}{\partial Z_i^{(2)}(T)} = \frac{\partial \text{DIV}}{\partial Y_i(T)} \frac{\partial Y_i(T)}{\partial Z_i^{(2)}(T)} \quad \text{OR} \quad \frac{\partial \text{DIV}}{\partial Z_i^{(2)}(T)} = \sum_j \frac{\partial \text{DIV}}{\partial Y_j(T)} \frac{\partial Y_j(T)}{\partial Z_j^{(2)}(T)}$$

$\underbrace{\hspace{10em}}$ scalar activation. $\underbrace{\hspace{10em}}$ vector activation.

$\sim \sim$

$$\begin{aligned} \nabla_{h(T)} \text{DIV} &= \nabla_{Z^{(2)}(T)} \text{DIV} \nabla_{h(T)} Z^{(2)}(T) \quad \nabla_{h(T)} \text{DIV} - \text{row vector} \\ &= \nabla_{Z^{(2)}(T)} \text{DIV} W^{(2)} \end{aligned}$$

$$[\quad] \quad [\quad]$$

Scalar form: $\frac{\partial \text{DIV}}{\partial h_i(T)} = \sum_j \frac{\partial \text{DIV}}{\partial Z_j^{(2)}(T)} \frac{\partial Z_j^{(2)}(T)}{\partial h_i(T)} = \sum_j w_{ij}^{(2)} \frac{\partial \text{DIV}}{\partial Z_j^{(2)}(T)}$

$\sim \sim$

$$\begin{aligned} \nabla_{W^{(2)}} \text{DIV} &= \nabla_{W^{(2)}} Z^{(2)}(T) \nabla_{Z^{(2)}(T)} \text{DIV} \\ &= h(T) \nabla_{Z^{(2)}(T)} \text{DIV} \end{aligned}$$

Scalar: $\frac{\partial \text{DIV}}{\partial w_{ij}^{(2)}} = \frac{\partial \text{DIV}}{\partial Z_j^{(2)}(T)} h_i(T)$

$$(*) \nabla_{z^{(1)}(T)} \text{DIV} = \nabla_{h(T)} \text{DIV} \nabla_{z^{(1)}(T)} h(T)$$

scalar form:-

$$\frac{d\text{DIV}}{dz_i^{(1)}(T)} = \frac{d\text{DIV}}{dh_i(T)} \frac{dh_i(T)}{dz_i^{(1)}(T)}$$

$$(*) \nabla_{w^{(1)}} \text{DIV} = \nabla_{w^{(1)}} z^{(1)}(T) \nabla_{z^{(1)}(T)} \text{DIV}$$

$$= X(T) \nabla_{z^{(1)}(T)} \text{DIV}$$

scat: $\frac{d\text{DIV}}{dw_{ij}^{(1)}} = \frac{d\text{DIV}}{z_j^{(1)}(T)} x_i(T)$

$$(*) \nabla_{w^{(1,1)}} \text{DIV} = \nabla_{w^{(1,1)}} z^{(1)}(T) \nabla_{z^{(1)}(T)} \text{DIV} \quad (\text{gradient and recurrent weight})$$

$$= h(T-1) \nabla_{z^{(1)}(T)} \text{DIV}$$

scat:- $\frac{d\text{DIV}}{dw_{ij}^{(1,1)}} = \frac{d\text{DIV}}{dz_j^{(1)}(T)} n_i(T-1)$

(*) (Q1): At this stage; we have only gone back one time step

- difference between previous presentations of conventional (MLP) backprop

CNN backprop:-

- Have to backprop all the way (to be clarified)

(*) omitting scalar form from hereon:-

$$(*) \nabla_{z^{(2)}(T-1)} \text{DIV} = \nabla_{y(T-1)} \text{DIV} \nabla_{z^{(2)}(T-1)} y(T-1)$$

\rightarrow S.2020 typo
Z

(analog to earlier)
scalar form for both vector act/
convent.)

$$\nabla_{h(T-1)} \text{DIV} = \nabla_{z^{(2)}(T-1)} \text{DIV} \underbrace{\nabla_{h(T-1)} z^{(2)}(T-1)} + \nabla_{z^{(1)}(T)} \text{DIV} \underbrace{\nabla_{h(T-1)} z^{(1)}(T)}$$

$$= \nabla_{z^{(2)}(T-1)} \text{DIV} w^{(2)} + \nabla_{z^{(1)}(T)} \text{DIV} w^{(1,1)}$$

scalar form: $\frac{d\text{DIV}}{dh_i(T-1)} = \sum_j w_{ij}^{(2)} \frac{d\text{DIV}}{dz_j^{(2)}(T-1)} + \sum_j w_{ij}^{(1,1)} \frac{d\text{DIV}}{dz_j^{(1)}(T)}$

⑦ note this has
the above functional form :- examine route of influence
from $h(T-1)$ to DIV (through 2 distinct routes)

⑧ view with slides ① - v. important

⑨ :- $\nabla_{w^{(2)}} \text{DIV} + \nabla_{w^{(2)}} z^{(2)}(T-1) \nabla_{z^{(2)}(T-1)} \text{DIV}$

\Rightarrow $\overset{(PK)}{\nabla_{w^{(2)}} \text{DIV}} + n(T-1) \nabla_{z^{(2)}(T-1)} \text{DIV}$

scalar form: $\frac{d\text{DIV}}{dw_{ij}^{(2)}} + \frac{d\text{DIV}}{dz_j^{(2)}(T-1)} n_i(T-1)$

⑩ computing $\nabla_{w^{(2)}} \text{DIV}$ requires summation / update
as the weight is identical to what was computed previously
 $w^{(2)}$ at time step T

$$\left(\nabla_{w^{(2)}} \text{DIV} = h(T) \nabla_{z^{(2)}(T)} \text{DIV} \right)$$

⑪ this is a significant change.

- similar to before :- $\nabla_{z^{(1)}(T-1)} \text{DIV} = \nabla_{h(T-1)} \text{DIV} \nabla_{z^{(1)}(T-1)} h(T-1)$
(some logic)

scalar form:-

$$\frac{dDIV}{dz_i^{(l)}(T-1)} = \frac{dDIV}{dh_i(T-1)} \frac{dh_i(T-1)}{dz_i^{(l)}(T-1)}$$

$\textcircled{1}$: Similar principle of 'inputting' applies as $\textcircled{2}$, :-

$$(*) \nabla_{W^{(l)}} DIV += \nabla_{W^{(l)}} z^{(l)}(T-1) \nabla_{z^{(l)}(T-1)} DIV$$

$$\Rightarrow \nabla_{W^{(l)}} DIV += X(T-1) \nabla_{z^{(l)}(T-1)} DIV$$

$$\frac{dDIV}{\nabla_{ij}^{(l)}} += \frac{dDIV}{dz_j^{(l)}(T-1)} x_i(T-1)$$

(*) And also for recurrent weights:-

$$\begin{aligned} \cdot \nabla_{W^{(l+1)}} DIV &+= \nabla_{W^{(l+1)}} z^{(l)}(T-1) \nabla_{z^{(l)}(T-1)} DIV \\ \Rightarrow (*) \nabla_{W^{(l+1)}} DIV &+= n(T-2) \nabla_{z^{(l)}(T-1)} DIV \end{aligned}$$

scalar form:-

$$\frac{dDIV}{dw_{ij}^{(l+1)}} += \frac{dDIV}{dz_j^{(l)}(T-1)} n_i(T-2)$$

(...)
Keep iterating using all principles above until reach initial state $h-1$

$\textcircled{3}$

$$\begin{aligned} \nabla_{h_{-1}} DIV &= \nabla_{z^{(l)}(0)} DIV \nabla_{h_{-1}} z^{(l)}(0) && h-1 \text{ is a param. equiv.} \\ &= \nabla_{z^{(l)}(0)} DIV W^{(l+1)} \end{aligned}$$

$$\frac{dDIV}{dh_i(-1)} = \sum_j w_{ij}^{(l+1)} \frac{dDIV}{dz_j^{(l)}(0)}$$

BRNN pseudocode

(A10)-review

- why are there many $t=$ \rightarrow resolve
- which ones are incrementing; and which aren't?
- ① A review why ~~$n \rightarrow$~~ n derivative of n is incremented (2)(?)

(*) Extensions to RNN: Bidirectional RNN (BRNNs)

(*) BRNN

- RNN with forward and backward recursion (2)(?)
- explicitly models the fact that just as the future can be predicted from the past; past can be deduced from future (???)

Read paper

- use future output to make inferences

(*) BRNN

- forward net processes input from $t=0 \dots t=T$
- backward $\text{---} \quad \text{---} \quad \text{---} \quad t=T, \dots, t=0$

(*) BRNN: Processing input

(?) represented

- forward pass } through BRNN produces ² hidden state _n rep of input.
- backward pass }

(!) Note this
is not backward!

(*) Output uses both 'hidden state representations'

(*) BRNN - pseudocode

(A11)-review

- (6m): forward pass
back pass

- (6): affine combination of output is weighted combination of forward pass with own set of weights, and backward pass (with own set of weights) and bias
- uses hidden state rep from forward and backward pass

(7) Bidirectional - simplified code

(8) A12 : clever flipping to implement - review.

(8) Backprop in BRNNs

- sees sides
forward
(6): propagate inputs through forward network, backward network.
" to get final output; compute divergences.

(9) Backprop:-

- 1) Forward net - ignore backward network; perform BPTT to get derivatives for the forward network; update parameters.
- 2) Backward net - ignore forward network; perform BPTT to get derivatives for back network; update params.

(X) Backprop for BRNNs

2 step:-

- i) BPTT subroutine
- ii) Param updates

(A13) - review pseudocode

