

Josh Starmer / Stat Quest

Neural Networks Series

①

15 video to cover

(More videos but outside of scope for
NNs, CNNs, RNNs, LSTM)

①	the main ideas of NNs	19
2	the chain rule	19
3	Gradient Descent	24
4	Backprop main ideas	18
5	Backprop details - opt 3 params Simult	19
6	Backprop details II - chain rule	14
7	ReLU	9
8	Multi-inputs & Outputs	14
9	Argmax & Softmax	15
10	Softmax derivative	8
11	Cross Entropy	10
12	Cross Entropy Derivatives & Backprop	23
13	Convolutional NNs	16
14.	Recurrent NNs	17
15	long short-term memory	21
		246
		4.1 hours

(2)

The main ideas of NNs

can fit to non-linear distributions

Activation function is what allows this
Softplus, ReLU, Sigmoid

Input \rightarrow hidden \rightarrow output nodes

Linear calc of input * weight + bias = x-axis coord

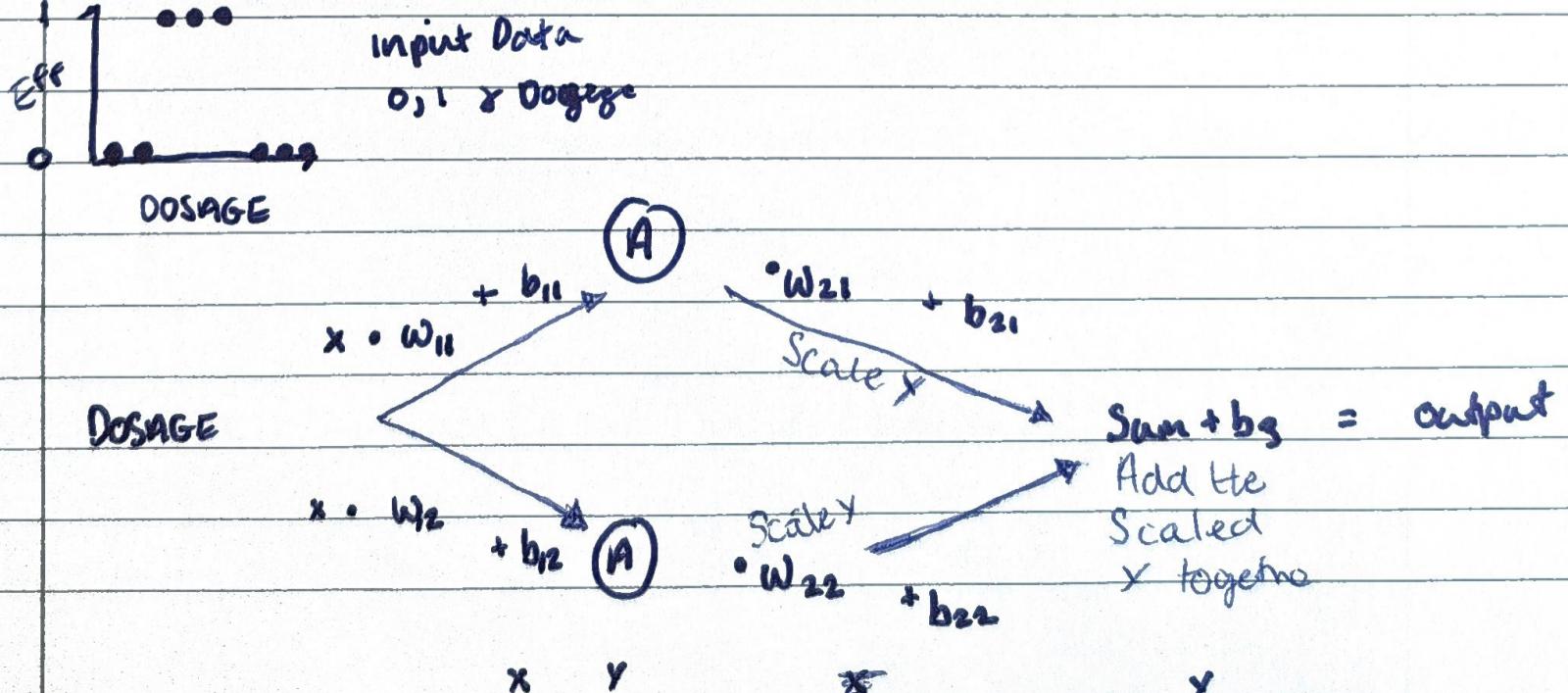
Put x into activation function = y-axis coord

the y-axis value is what is then carried into the next layer & the process repeats

Plot y-axis curves to see the space the area of the network is working in

Add together the y-axis curves from different routes to create a new curve
(And minus the bias).

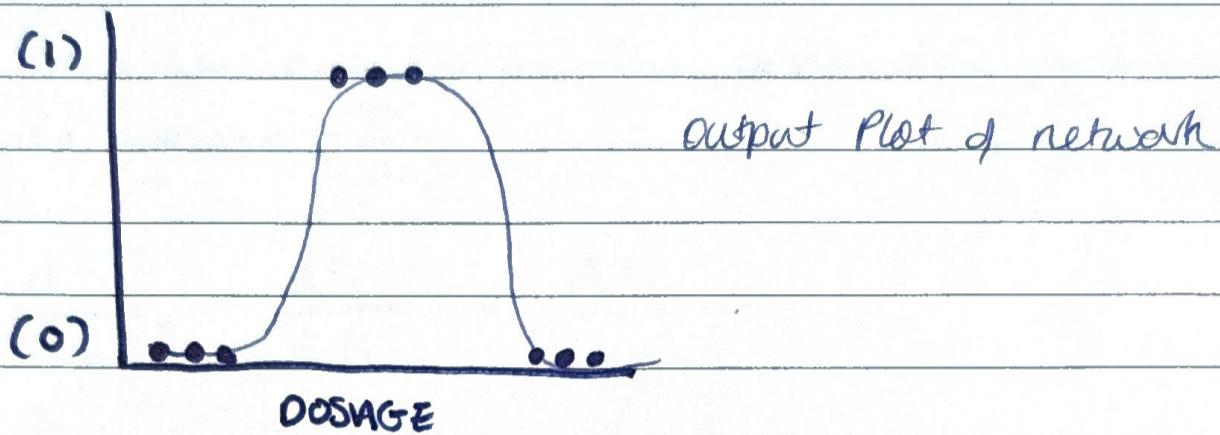
This is the output of network that fits the data



(3)

two options now:

- Plug new values into network to get output
- Plot the graph the network creates & read from it



weights, bias, nodes

② Chain Rule

Assumes basic understanding of derivatives

Chain rule is used for intertwined relationships

Weight & height \rightarrow height & shoe size

using height to predict shoe size

1 unit increase in weight = 2 unit increase in height

$$\frac{dH}{dw} = \frac{2}{1} = 2$$

$$\text{So Height} = \frac{dH}{dw} \times w = 2 \times w$$

1 unit increase in H = $\frac{1}{4}$ increase in w.

$$\frac{1}{4} \frac{DSS}{DH} = \frac{1}{4}$$

$$S = \frac{DSS}{DH} \times H = \frac{1}{4} * H$$

(4)

the two can then be interwoven

$$S = \frac{dS}{dH} * H + \frac{dH}{dw} * w$$

now we can take the derivative of shorsize w/ resp
to weight $\frac{dS}{dw}$

$$\frac{dS}{dw} = \frac{dS}{dH} * \frac{dH}{dw}$$

(3) Gradient Descent

A lot of ML is about optimization

$$\text{Pred} = \text{Intercept} + \text{Slope} * \text{weight}$$

to start, pick an slope and plug in a slope

then pick an intercept @ random. We are going to use GD to find the correct intercept

provides a guess for GD to improve from

- ① Calculate loss functions of ~~your~~ data points
 - Against the prediction line

$$\text{Residual} = \text{Actual} - \text{estimate}$$

$$\text{Sum of squared residuals} = \text{loss function}$$

Plot graph of loss func w/ diff intercepts

this demonstrates that there is a curve
at some point exists an optimal intercept
to achieve lowest loss function

But How to move to this point?

- GD

take the derivative of the loss function
w/ respect to the intercept

e.g loss function =

$$(1.4 - (\text{int} + \text{slope} * \text{pred}))^2 +$$

~~

~~

(use chain rule)

~~$\frac{d}{d\text{int}}$~~ Sum of squared residual
~~func~~

GD uses the derivative to find the optimal point

Plug the intercept (o) into the derivative
to get the slope of the curve at that point

note when close to the optimal point the
slope will be close to 0

higher slope when further away

this provide the size of the "Step" for
GD to move the intercept for a better
value

GD takes the slope, e.g. 5.7, & applies a "learning rate" e.g. 0.1

$$\text{Step size} = -5.7 * 0.1 = -0.57$$

$$\begin{aligned}\text{new Intercept} &= \text{old} - \text{step} \\ 0 &- (-0.57) = 0.57\end{aligned}$$

then iterate to learn again

How many steps to take? ↗ e.g.

Stop when step size = 0.0001

How to estimate slope & intercept?

Same as before but this take the Deriv of the loss function w/ resp to intercept and also the slope

$$\frac{d}{d \text{Inter}} \quad ? \quad \frac{d}{d \text{Slope}}$$

having to derivatives of the same function is known as a gradient

use gradient to descend to lowest point in loss f
- Hence the name

Pick random intercept (0) & slope (.)

Plug these into the 2 Derivatives & get out 2 Slopes

(7)

turn these slopes into steps w/ LR

$$\text{Step size}_I = \text{slope} * \text{LR}$$

$$\text{Step size}_S = \text{slope} * \text{LR}$$

use the slope size to update both I & S

Repeat steps until step size \rightarrow ^{small} ~~really small~~

Steps:

- ① take derivative of the prep function for each param inside of it
- ② pick random starting values for the parameters
- ③ plug the params into derivatives
- ④ calculate the step sizes
- ⑤ calculate/update the parameters

$$\text{New } P = \text{old } P - \text{step size}$$

- ⑥ keep repeating 5-5 until some small step size

④ Backprop Main ideas

- ① use chain rule to calc derivs
- ② use gradient desc to update

Backprop starts w/ the last param

Recall that activation funcs output y-axis coords which we can plot

each netpath in the net creates its own y-axis curve

- A curve is obtained by plugging in several inputs, e.g. 0-1 in 0-1

A curves from DIFF routes together to get final fitting curve

+ Add final bias

- what this example wants to optimize

Set at random value 0

Calculate loss function against network

Plot loss against bias (0)

- increment Intercept to demonstrate how intercept can change loss func

What to use gradient descent to make this movement

find the derivative of the loss function
w.r.t respect to intercept

$$\frac{dL_F}{db_3} = \frac{dL_F}{d\text{Pred}} * \frac{d\text{Pred}}{db_3}$$

Chain
rule to
get $\frac{d\text{Pred}}{db_3}$

After follow known gradient desc rules

⑤ Backprop Details Pt 1 - Opt 3 Paras Simultaneously

How to get ~~all~~^{more} other optimal params
- not just the last bias

Calc b_3, w_3, w_4 @ once

Start init w_3, w_4 w/ random vals

init B_3 as 0
- not it is common to init B
intercept as 0

Run network w/ input & get output
using initd params

Calculate the loss function

use same steps as before to GD b_3

Important point here is that even though
we are optimizing further down the tree
the calc's are the same

how to calc derivatives for w_3, w_4 ?

breakdown how we get prediction rate

$$\text{Pred} = \text{Activ}_1 * w_3 + \text{Activ}_2 * w_4 + b_3$$

↳ This value is used to create the LF

And writing it this way demonstrates how to link the LF to w_3, w_4

Can now use chain rule to:

$$\frac{d\text{SSR}}{dw_3}$$

$$\frac{d\text{SSR}}{dw_4}$$

$$\frac{d\text{SSR}}{dw_3} = \frac{d\text{SSR}}{d\text{Pred}} * \frac{d\text{Pred}}{dw_3} \quad \begin{array}{l} \rightarrow \text{easy to} \\ \text{Derive -} \\ \text{Pred Derived} \\ \text{by } w_3 \end{array}$$

$$\frac{d\text{SSR}}{dw_4} = \frac{d\text{SSR}}{d\text{Pred}} * \frac{d\text{Pred}}{dw_4}$$



note uses the same calc as each other but
Also the same from the b_3 calc

Now we have 3 derivatives to plug into
& gradient descend

$$\frac{d\text{SF}}{dw_3} = \sum_i -2 \times (\text{obs}_i - \text{Pred}_i) \times \text{Activ}_1 / \star$$

Expand, Plug observed, Plug Observed, Plug Value

Output is a slope value
 → Step size

Repeat for w_4, b_3

Apply GD update process

⑥ Solving all parameters

Optimize all weights & bias

Want to get to w_1 :

$$\frac{d\text{err}^L}{dw_1} = \frac{d\text{err}}{d\text{pred}} \cdot \frac{d\text{pred}}{dy_i} \cdot \frac{dy_i}{dx_i} \cdot \frac{dx_i}{dw_1}$$

\downarrow \downarrow \downarrow
 $y_3, w_3 +$ activfunct
 $y_2, w_2 +$ $\log(1 + e^x)$
 b_3 Input * $w_1 +$
 b_1

Follows the connections back

Note all derivatives are the same structure regardless of the number/function seen

This is why Backprop can be easily coded

Plug all derive together to get a new function of the Deriva $\frac{d\text{err}}{dw_1}$

Repeat process for b_1

Note the steps repeat most of the same steps as w_1

Same steps for w_2 & b_2

Then use GD to opt all params together

Start by init weights & bias

- weight = ~~rand~~ rand norm dist (many ways to do this)
- bias = 0's

Expand Derivative summations and plug values
to get predicted values (from network)

Solve each derivative & turn into step size
Step size = Derivative * learning rate

& update params

$$\text{new} = \text{old} - \text{Step Size}$$

Repeat until criteria met
(Step size $< \alpha$)

⑦ Relu in action

Activation function

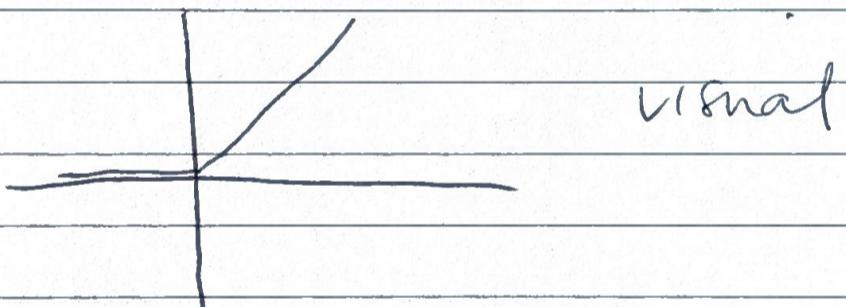
Note can also use activation functions at the final output too

$$y = f(x)$$

$$\text{Relu} = \max(0, n)$$

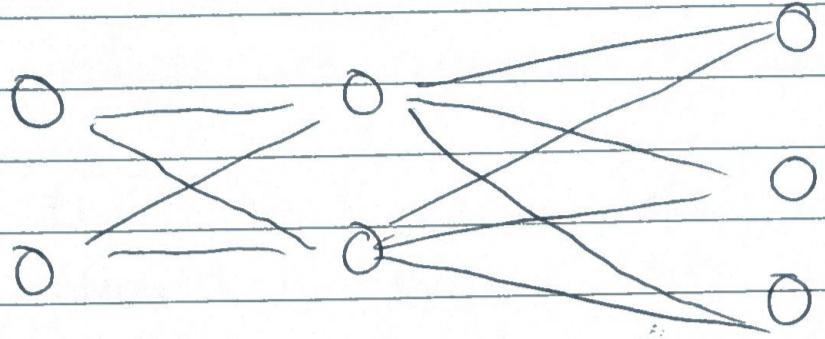
if less than 0 ~~then~~ then $\text{Relu} = y = 0$

if > 0 then $y = x$



(1x)

③ Multiple ~~step~~ inputs & outputs



All interconnected w/ weights

for the output active func common to
convert values to $0 - 1$

this 1 is more likely to be the correct value

Select the highest (Argmax or Softmax)

(9) Argmax & Softmax

labels are denoted by a \mathbf{z}

thus nn outputs for categorical data
 Should be below 1 hence normalization
 activ func between 0 & 1
 or use Argmax

Argmax

- Sets all outputs to 0 or 1
- targets 1, rest = 0
- then can be used w/ loss function
 $-(\text{obs} - \text{Pred})$
- this can't be used to optimize a nn
 because the deriv of 0 is 0
- no info \rightarrow no grad Desc step

Softmax

- Argmax for output, Softmax for training
 $v_i \leftarrow \text{one cat}$

$$v_1 + v_2 + v_3 \leftarrow \text{all cats}$$

- Always ≤ 1 & sums to 1 over all cats
- Can be interp as "probs"
- Don't trust as true probs

$$\frac{e_i}{\sum e_j}$$

App

- Deriv of Softmax is not always 0 &
 hence suitable for gDesc

(10) Soft-Max Derivative

Derivative of category X

with respect to the

Derivative of the raw output value of Cat X

$$\frac{d "P_{\text{Cat } X}"}{d \text{ Raw Setosa}} = "P_{\text{Cat } X}" \times (1 - "P_{\text{Cat } X}")$$

"P Cat X" is the probability from the SF
 ↳ Deriv of this use Quotient rule to calc
 (complicated)

raw output is just the pre-SF value

Important = each SF outvalue relies
 on all of the categories

CAT A

CAT A + CAT B + CAT C

So need to take the deriv w/ respect to ea

$$\frac{d \text{CAT A}}{d \text{raw A}}$$

$$\frac{d \text{CAT B}}{d \text{raw B}}$$

$$\frac{d \text{CAT C}}{d \text{raw C}}$$

draw A

draw B

draw C

$P(\text{CAT A}) \times P(\text{CAT B})$ easier than own Cat

(17)

(11) Cross Entropy Loss function

w/ one output we often use sum [1/reds]

w/ 2+ we want to use Argmax as the output activation

↳ has a bad derivative so can't use for Backprop

so we use softmax which outputs values as probabilities, 0-1

when outputs are bounded between 0-1 we use cross entropy cf

Raw output \rightarrow softmax Probs

$$\text{cross ent} = -\log(\text{softmax})$$

note $-\log(x)$ isn't the standard form of cross entropy

NNs only need a simplified version of it

obtained by expanding & summing

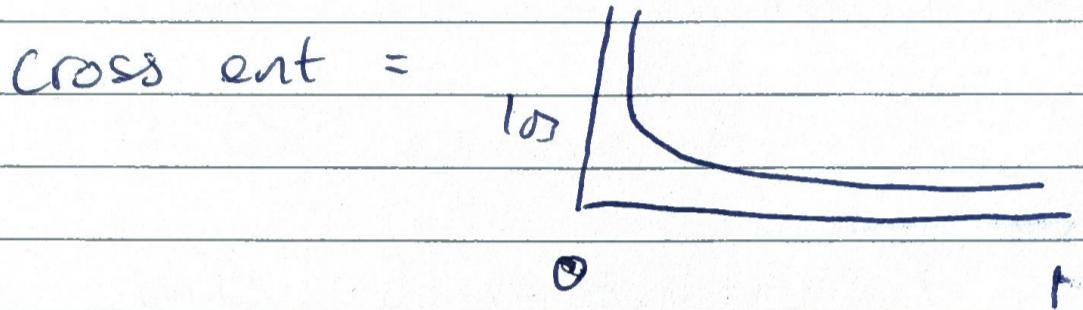
$$-\sum_n \text{observed}_n \times \log(\text{predicted}) = \text{full}$$

$$\text{total error} = \sum_n (-\log(\text{sf}))$$

From the indiv cross ent values we are meant to calc the sum of square resid

why not just calc w/ softmax probs?

softmax probs = $0, \phi$



As n cross Ent get more bad, it gets an exponentially higher value

this provides more info to learn / penalise

the slope of the derivative is larger @ bad preds

= larger step size to learn

(12) cross ENT & Backprop

raw \rightarrow softmax Prob \rightarrow entropy value

cross ent depends of max

but softmax has a different calc
depending on the observed value

~~cat(n)~~ cat r, n, m

thus different derivs to be calced

$$CE = -\log(\text{Pred Cat } n)$$

$$\text{Pred Cat } n = \text{softmax}(\text{cat } a, \text{cat } b, \text{cat } c)$$

Linking cross ent to paras

↳ find the link via softmax category

Derive uses chain rule

$$\frac{d \text{Cat}}{db_3} = \frac{d \text{Cat}}{d \text{Prob}} \times \frac{d \text{Prob}}{d \text{Raw}} \times \frac{d \text{Raw}}{d b_3}$$

Plug in & simplify

When obs cat is used to calc
the deriv then = prob - 1

for other cats still need to follow
the route for the correct cat

this means taking the derive via the
softmax

$$\frac{d \text{cat 2}}{d b_3} = \frac{d \text{cat 2}}{d \text{Prob C2}} \cdot \frac{d \text{Prob C2}}{d \text{raw C1}} \cdot \frac{d \text{RCI}}{d b_3}$$

Plug & Simple

Result = Prob cat 1

How to optimize?

Only plug in the correct derivative route
for each data entry w/ respect
to param

Derive of cross ent = slope

Step size = slope * learn rate

GD iterate

(13) convolutional NNs

Mainly used for image classification

Convert matrix of image into vector

This vector then proxies as the input layer

each pixel = input neuron

Plain NNs are poor for images because they do not scale well

each pixel / neuron has a weight to the next layer

big / detailed images = many pixels = many weights
= too comp heavy

For large & complicated images
→ CNNs

CNNs apply a filter to the image which is 3x3 pixels (kernel)

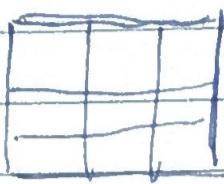
0	0	1
0	1	0
1	0	0

looks like this but is estab via backprop (?)

this kernel is then applied iteratively to the image

the result is a small image as

kernel



= n

9 pixels to

, i.e. 1 Pixel

9 Pixels overlaid w/ kernel & Dot product
(each Pixel times & added together)
→ Sum of Products

~~After~~ After applying the kernel add a bias term

Next slide the kernel over 1 pixel & repeat

- ↳ • The slide of 1 pixel can vary
- Some methods slide more to reduce the computation cost

why is compressing the image okay?

- Areas of a image are often correlated so info is not really lost

After kernel conversion an activation func is applied. (ReLU)

~~And then apply another kernel~~
method called Max Pooling

this takes a kernel, using of 4, and selects the highest value

usually this method has no overlap & slides to the next empty slot

Filter kernel \rightarrow ReLU \rightarrow Max Pool Kernel
 9 to 1 > 0 4 to 1

\downarrow
 the best filter match the
 input cluster

\downarrow
 Pooling select the
 most relevant /
 strongest value
 from the filter

the exists an alt method called mean pooling

finally the max pooled image \rightarrow be
 converted into a vector \rightarrow used
 as an input to a standard NN

this is fine as it is now essentially a
 small image

Output of NN is the classification of
 the data e.g. naughts or crosses

Image \rightarrow Filter to Feature Map \rightarrow ReLU \rightarrow Max Pool \rightarrow Input nodes

6x6 images is compressed down to 4 numbers

the filter kernel is the same as look @ a
 cluster of pixels rather then just 1

⑭ Recurrent Neural Networks (RNNs)

RNNs are a stand alone thing but they also feed into

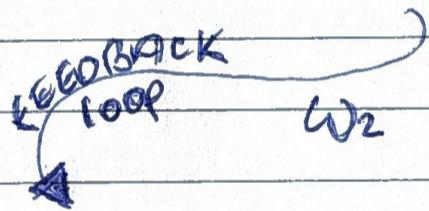
- long short term memory networks
- transformers

Varied input lengths?

RNNs have feedback loops

RNN looks like it has a single input value but the feedback loop allows for sequential or time series data

yesterday $\rightarrow w_1 \rightarrow b_1 \rightarrow$ Activ Y_1
input



today $\rightarrow w_1 \rightarrow b_1 \rightarrow$ Activ Y_2
input
Sum

\downarrow
 $w_3 \rightarrow b_2 \rightarrow$ Pred
Tommer

feedback loop allows additional layers to be threaded in a sequential data

the weights & biases w_1, w_2 are shared between layers

7

basic RNNs are not used often, why?

- the more we unroll a RNN the harder it is to train
 - ↳ known as vanishing/exploding gradient
- Descent problem

note W_2 is the feedback weight

7

if the W_2 is > 1 then the input value keeps getting multiplied by factor n
= A huge number

huge number would find its way into the gradient Descent & step sizes wouldn't be incremental

Solution could be to limit W_2 to $\{1$
but this leads to vanishing gradient

7

= steps are too small

LSTM NN are a solution & extension to this

D

(15) Long Short-term Memory (LSTM)

Solution to RNNs & Precursor to transformers

recap Vanilla RNN = feedback loop for sequence
 = explode or vanish gradient

RNNs = hard to train

LSTM designed to avoid

has two prediction loop routes:

- Feedback loop is used for short-term
- weights direct to output for long term

LSTM is a much more complicated unit

uses sigmoid & Tanh H activ funcs

Sigmoid = takes x-axis coord & creates y-axis between 0-1

tanh = x to $y = -1$ to 1

longterm has a lack of weight & bias

Short-term updates the long-term memory

- ST uses Sigmoid activ to be applied to LT
- So update can be thought as % being remembered

next unit of ST = Potential

- uses tanh to create a value -1 to 1
- & Sigmoid to % the potential value
- which then gets added to the LT

Updating the long term memory
- Sigmoid

Init input twice for LT & ST

for sequences, only the ST is inputted

the NN runs in "units" updating the LT
& ST as it goes and inputting the seq
ST data