

24/0

ML - Week - Advanced NN's

- Convolutional NNs (CNN)
- Recurrent NN (RNN)

CNNs

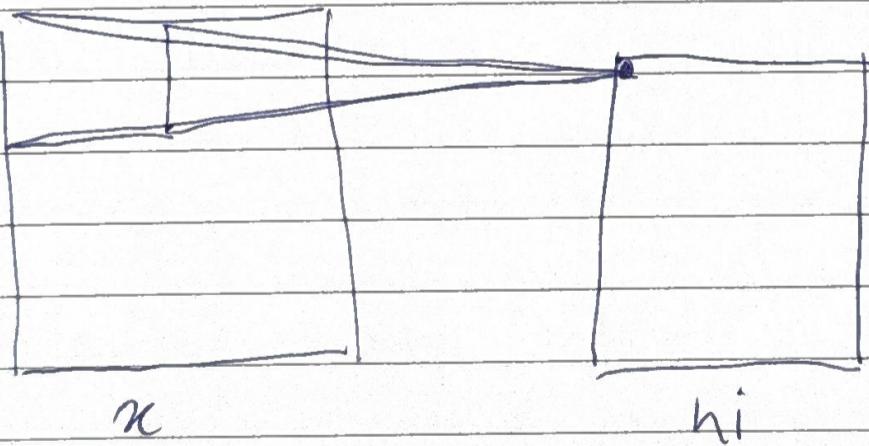
Continues the general concept of layers

Sections of the layers predict sections in the next layer

CNN created for grid structure data

learns ~~spat~~ spatial relationships

e.g. image \rightarrow neighbouring pixels are more related than distant pixels



Square \rightarrow Point

Compared to MLP layer:

MLP = unit \rightarrow unit in other layers

Convolution small weight matrix re-applied

- Convolutional has a neighbour structure
- MLP has unit to unit structure

Colour has 3 channels \rightarrow RGB

Image = matrix array \rightarrow rows = 3-D matrix

With Matrix = the rectangle

~~As~~ input = width matrix, output = single pixel

With Matrix = height, width, channel (color, shade etc)
 \rightarrow Should be significantly smaller than input
 data (image) h, w, c

Weights are called a kernel in CNN layers

~~each output Does each Pixel have a unique kernel??~~

Stride = how much the width matrix moves

Padding \rightarrow Allows better preds @ edge of image

Input data example: $w \times h \times c = 7 \times 7 \times 3$

each channel has a weight matrix

Activation funcs & CNNs:

- ReLU common
- Alternate layers between convolution & "Pooling"
- Pooling = aggregations to data
- lower size of image

Pooling allows invariance to minor transforms

MLP all dimensions are interconnected

~~Each~~ CNN uses smaller weights

- Week 8 lab =
- UK Met data
 - Applies convo to just 1 dimension

Question: is activation (relu) applied to pixels?

the use of CNNs is very expensive and reqs GPUs

2D CNN is the standard for image data

→ 3D for video where 3rd is time

Summary

Primary feature of weight sharing of CNNs

weight sharing reduces overfitting

hyperparams are kernel, stride, padding & channel

Recurrent Neural Networks (RNN)

- MCP has no connections within the same layer
- only between layers

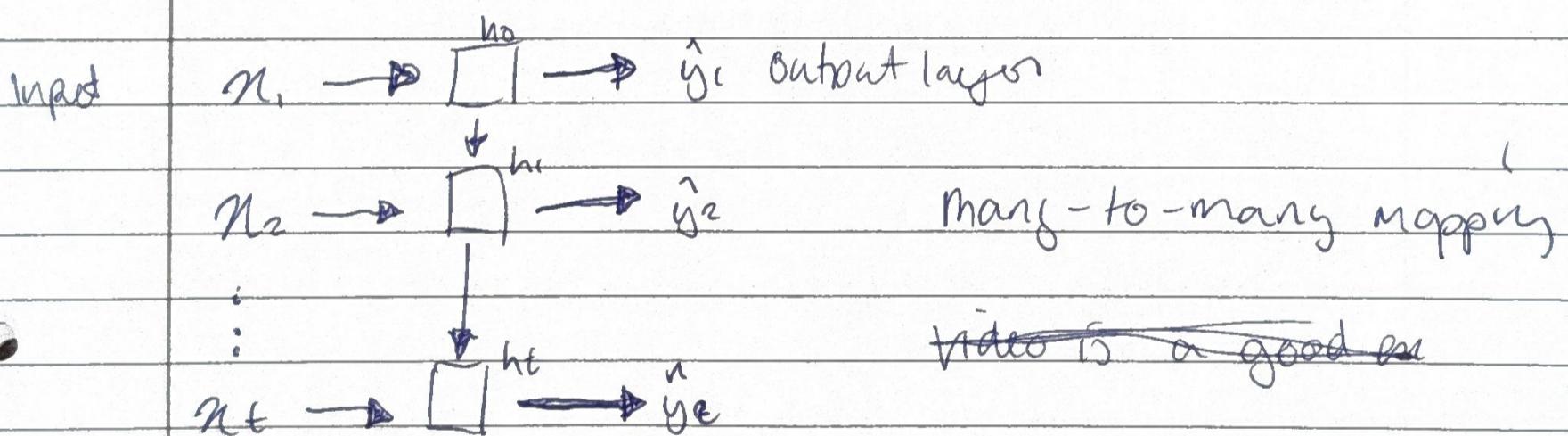
RNNs have connections between layers

this connection is feedforward

Designed for sequential data

Recurrence function is always the same

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



A common example is video due to time aspect

Another is prediction of the next sequence of words given an input sequence of words

RNN can be reduced to many-to-one in order to just predict the \hat{y}_T

or one-to-many, e.g. image ~~to~~ ^{to} text / summary

3

Problem w/ RNN:

vanishing gradient

Recall stochastic gradient descent

in MLP, gradient is backpropagated ~~to the~~ from the loss func to the input \Rightarrow update

RNN w/ time data = Backprop through time

Problem is when you go further enough back in time, the gradient becomes very small \Rightarrow almost 0

this means the model stops learning

long short-term memory (LSTM)

As a solution to vanishing gradients

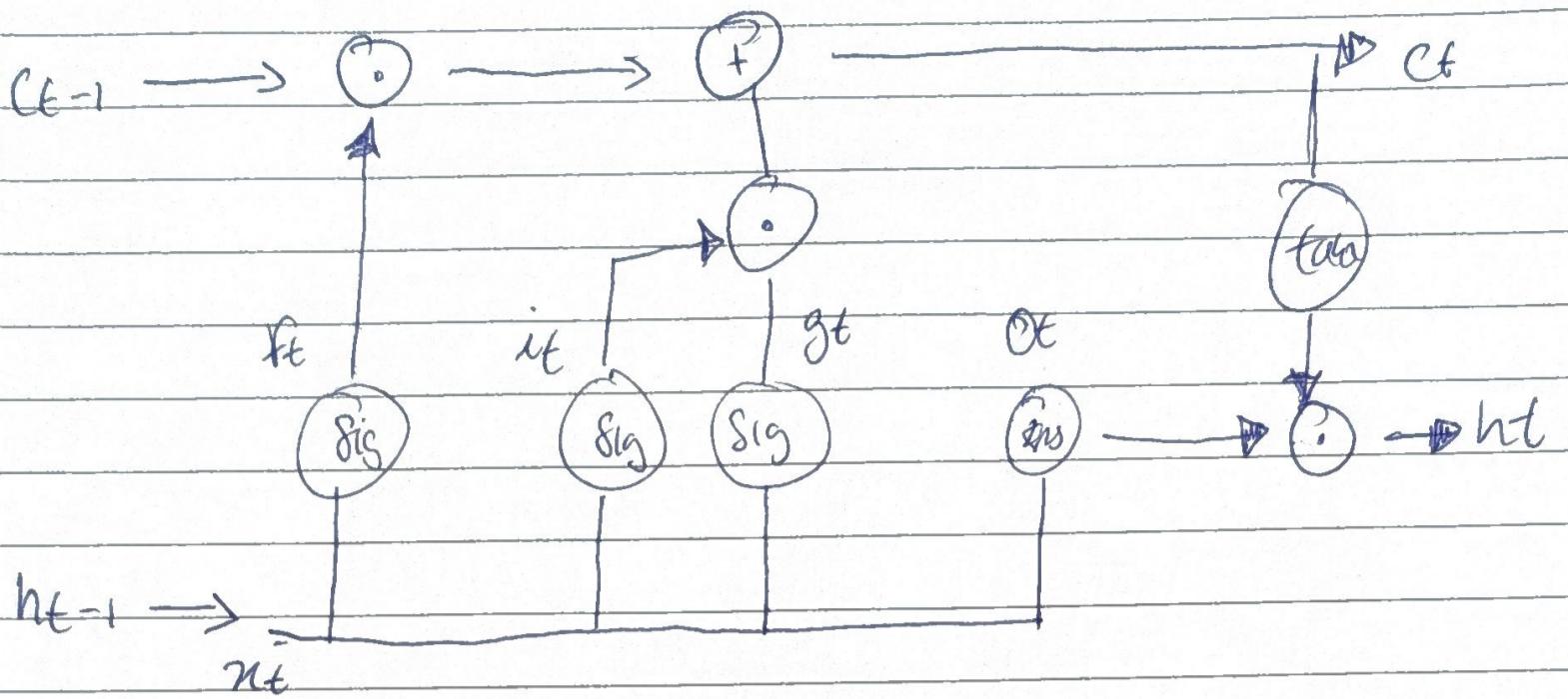
$h_t =$

before: $f(h_{t-1}, u_t, w)$

After: $h_t = f(c_t, o_t)$

uses gates to control flow through time

$$c_t = f^*(c_{t-1}, f_e) + f^*(i_t, g_t)$$



forget gate = $f_t = f_c(h_{t-1}, x_t, w_f)$

Block input gate = $i_t = f_i(h_{t-1}, x_t, w_i)$

Cell update = $g_t = f_c(h_{t-1}, x_t, w_c)$

output gate = $o_t = f_o(h_{t-1}, x_t, w_o)$

How does it does it correct vanishing grads?

- LSTM has blocks (usually 1 but can be multi)
- each Block has multi-cells
- cells in Block share gates

Error flow in LSTM:

- Backprop within cells & Blocks
- No Backprop through time

↳ i.e. Problem addressed