

①

# E2 Learn AI - Deep Learning 101 .

- ① Convolutional NNs
- ② RNNs
- ③ LSTM
- ④ Cross entropy
- ⑤ Stochastic Gradient Descent
- ⑥ Autoencoder

# CNNs

works very well for image recognition & pattern recognition

CNNs learn by extracting & learning spatial ~~fea~~ features of image data

Steps:

① convolutional layer (kernel)

② Pooling layer

③ fully connected layer

kernel = filter layer, elementwise, sliding, makes matrix smaller

Image  $\xrightarrow{\text{kernel}}$  feature map

Different types of kernel can highlight different types of lines in an image - (horiz, vert, curved, etc)

kernel size impacts size & detail of the feature map more

small kernel = bigger, more detailed feature map  
 $\rightarrow$  but more computationally expensive

$\text{Stride} = \text{Slide of kernel}$

bigger strides = less detail & smaller feature map

Convolutional layers introduce an activation func to handle non-linearity

→ ReLU most typical

### ► Pooling layer

Job is to reduce the size of the Feature Map

main = Max Pooling or Avg Pooling

use kernels w/ no stride

Max = Select max in kernel

Avg = Average of all values in kernel

### ► Dense / Connected layer

Creates final output based on extracted features

Conv → Pool → Conv → Pool → Flat → Dense

training of a conv involves changing the weights of the kernel in the Conv layer(s) to reduce error of output prob

Backprob & GD used

Over time/training the kernel is optimiz to distinguish between the local feat of  $x$  &  $z$

## DNNs

Process sequential data

↳ time series data

↳ continuous information

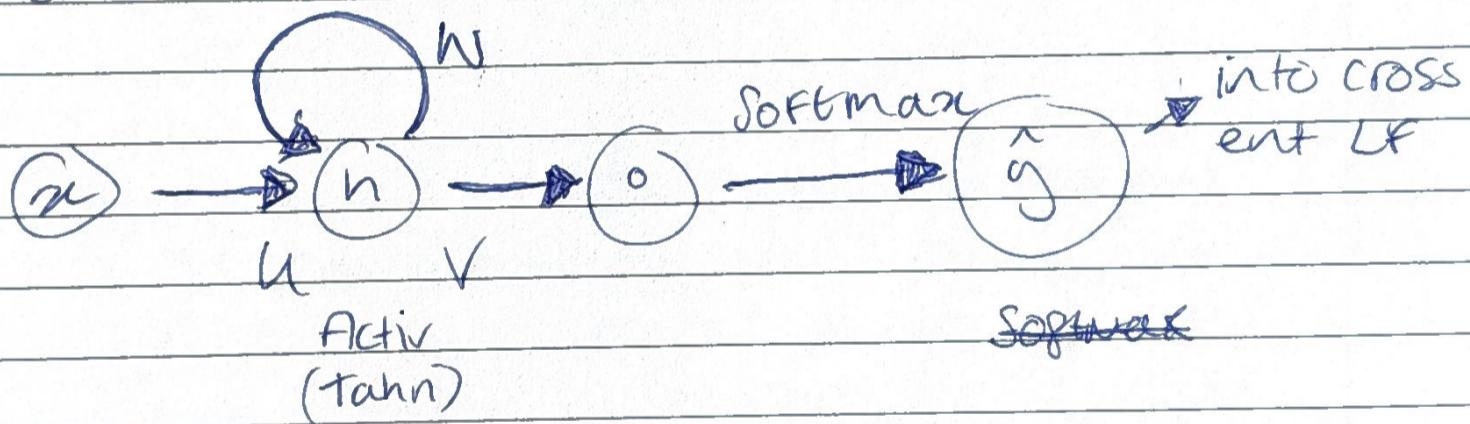
→ Stock prices

- Words/Sentences

if CNNs extract temporal spatial features from images

RNN extracts temporal features from time series data

RNN → LSTM → transformer



$h_1$  gets calc'd & passed into  $h_2$

$u, v, w$  are weights

## LSTM

RNN weakness = longterm dependency  
the longer the time, the more chain rule

\* Cells inner

- if  $L_1 = \text{diminishing}$  } early values also
- if  $T_1 = \text{explode}$  } get lost

to overcome = LSTM

how? = cell state

4 gates

- Forget gate
- Input gate
- Candidate gate
- Output gate

Forget gate = Decides what to forget

- takes input as prev state } concat
- At new/current state } matches
- Sigmoid ( $W \cdot \text{concat}$ ), 0-1

↳ meet cell state & updates element wise

Input gate

- same calc,  $\sigma(w_i x_i)$ ,  $w_i$  weights

Cand gate

- Work w/ input, elementwise
- tanh( $w_j x_j$ ) (-1, 1)

Add & update cell state

CST = long term memory

HST = short term (Hidden)

two flows of information

Note: video has good example of  
setting up & initing weights  
to calc the Read forward  
(qun)

HST is derived from the CST layer  
+ the output / readi gate

Output =  $z_t$  comes from only HST  
 $z_t = w \cdot H_{8t}$  (raw logit output)

$\xrightarrow{Softmax} \hat{y}_t$

CST & HST Shape do not need to equal  
the shape of output vector  
↳ the  $z_t$  weights will convert  
to the right size

5 sets of weights to Backprop & Grad Desc

- Forget G
- Input G
- Candi. G
- Output G
- $z_t$  (raw output)

for each section of the model separate  
out the weights calc & activation  
for back prop to take place

Softmax & cross ent are best pair as

$$\frac{\partial L}{\partial z_t} = \hat{y} - y$$

Note: video has good breakdown of calc  
the derivatives w/ mention of  
caveats to consider

## Cross-Entropy loss

frequency based loss func

$$H(p, q) = - \sum p(x) \log q(x)$$

~~definition~~

Aspects of cross ent:

- Information
- Expectation value
- Entropy

information (surprise)

Prob & Info are inversely related

$$P(x) \rightarrow 1/P(x)$$

In Info theory the actual formula is =

$$\bullet \log_2 \left( \frac{1}{P(x)} \right)$$

In info theory, information is a measure of surprise / unexpectedness expressed as an ~~abs~~ objective number

$$\log \left( \frac{1}{P(x)} \right)$$

Degree of Surprise

$$\text{or } \log(1) - \log(P(x))$$

$$0 - \log(P(x))$$

$$-\log(P(x))$$

## Expectation Value

Skill level \* probability = Expected Value  
Objective

EV considers the objective data & the daily (?) variables

EV = value Expected today

$$E(x) = \sum x_i \cdot p(x)$$

## Entropy

linked to EV & Surprise

take ev calc & replace w/ surprise

$$H(x) = \sum \log\left(\frac{1}{p(x)}\right) \cdot p(x)$$

"Expected degree of surprise"

why & when to loose?

→ consider football example where both  $E(x)$  is the same for both teams

then consider Entropy

- if prob is near 0 or 1 = predictable
- near 0.5 = unpredictable

More predictable  $\neq$  more chance of win

it can be an important datapoint in det  
the outcome of the match

Entropy is used to calc this type of  
predictability

More predictable = low entropy

high Ent = unexpected events occur  
frequently

## Cross entropy

$$\sum \log\left(\frac{1}{p(x)}\right) \cdot p(x) = \text{entrohs}$$

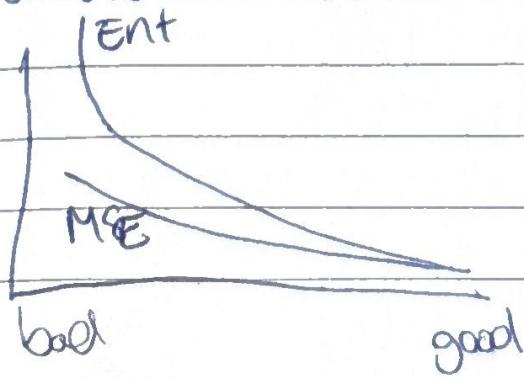
cross ent change P to q

$q(x)$  = Surprise of  $x$

Multiply Surprise of  $x$  by prob of  $x$

Why X-ent better than MSE?

MSE loss better than MSL



@ bad predict the gradients  
one much steeper

better for learning

# Stochastic Gradient Descent

3 types:

- Batch
- Stochastic
- mini-batch

NNs inherently have the local minima problem

→ no true solution but these 3 are methods explored

known as data scheduling for training

Batch = use all data points to calc  
the derives for the loss func  
for a single weight update

~~Stochastic~~ updates in <sup>manor</sup> stable ~~many~~ but  
is very inefficient if many data points

Stochastic = calc error for just 1 point  
& update the weight

- More efficient
- Very sensitive to each point  
↳ increase local minima problem

Batch is more stable & less prone to local minima

### mini Batch

uses set num of data points for each datapoint

- More stable than Stoch
- Mae effects that Batch

Note: has a code example