

# Convolutional Network

## Convolutions in 1D:

- Impulse response
- Multiple Kronecker Delta
- Gaussian Averaging

## Linear Local Agg.

$$f(x_i) = ax_{i-1} + bx_i + cx_{i+1}$$

## Activation Functions

$$y = \tanh(x)$$

$$y = \text{sigmoid}(x)$$

$$y = \text{ReLU}(x)$$

- Very Fast
- Gradient
- Piecewise linear interpolation

## Max vs Avg Pooling:

Max pooling is more sensitive to noise

## Why MLE won't work?

- Large dimensions
- No. of parameters goes to huge range
- Not easy to learn invariance.

Let's assume input is  $W_1 \times H_1 \times C$   
 Conv layer needs 4 hyperparameters:  
 - Number of filters **K**  
 - The filter size **F**  
 - The stride **S**  
 - The zero padding **P**  
 This will produce an output of  $W_2 \times H_2 \times K$   
 where:  
 -  $W_2 = (W_1 - F + 2P)/S + 1$   
 -  $H_2 = (H_1 - F + 2P)/S + 1$   
 Number of parameters:  $F^2CK$  and  $K$  biases

CS=SC

## Correlation vs Convolution

- Circulant Matrix
- Shift Invariant

## Inputs & Outputs

INPUT:  $B \times C_{in} \times H_{in} \times W_{in}$

- ↓ Samples
- Multi-channel processing

CONVOLUTION LAYER:  $C_{out} \times C_{in} \times K_H \times K_W$

- ↓ Bias Parameters

OUTPUT:  $B \times C_{out} \times H_{out} \times W_{out}$

$C_{out}$  filters of shape  $C_{in} \times K_H \times K_W$

## Loss Function:

$$\rightarrow L = \|y - p\|^2 \rightarrow L_2 \text{ loss}$$

→ Surface is convex

→ MLE

↳ log likelihood:  $\log(p(y|x, \theta))$

$$L = \log p(y|x, \theta) = \sum_i \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i w)^2}{2\sigma^2}} \right)$$

$$\max_w L = \max_w \sum_i \log c - \sum_i (y_i - x_i w)^2$$

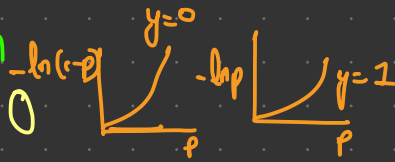
$$= \min_w \sum_i \|y_i - x_i w\|^2$$

↳ Same as MSE

→ Residual connections solve problem of vanishing gradients → Adds info

→ Parameter initialisation

→ Weights shouldn't be 0



→ 1x1 Convolution

INP:  $B \times C_{in} \times H \times W$

CON:  $C_{out} \times C_{in} \times 1 \times 1$

Converts  $C_{in}$  is converted to  $C_{out}$   
↳ MLP in spatial dimensions

→ Normalisation

$$x_j = \frac{x_i - \mu_i}{\sigma_i}$$

→ Batch norm:

① Normalise → Scale Shift

② Moving Average

→ Spatial Pooling

↳ Provides invariance

Binary Cross Entropy:

$$\rightarrow H = -p \log p \Rightarrow \text{entropy}$$

$$\rightarrow \mathcal{L}_{BCE} = -[y \ln p + (1-y) \ln (1-p)]$$

$$\hookrightarrow \{0, 1\} \quad p \in [0, 1]$$

When extended to K-classes:

$$\mathcal{L}_{CE} = -\sum_K y_K \log(p_K)$$

① CE maximises log-likelihood

② maximise log-likelihood = minimise  $-\text{ve}$  log-likelihood

K classes: SOFTMAX

$$\text{sigmoid}(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{1+e^z}$$

$$z \in \mathbb{R}^K, p = \text{softmax}(z)$$

$$p_i = \frac{e^{z_i}}{\sum_j e^{z_j}}, p_i \in [0, 1] \\ \sum_j p_j = 1$$

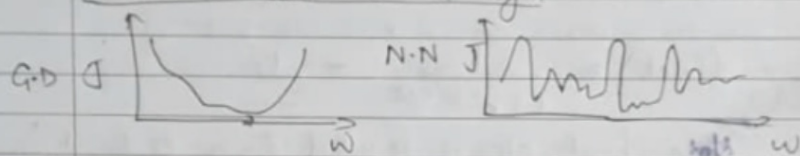
Examples:

① AlexNet ② VGG-Net

③ GoogLeNet ④ ResNet

# 6.8 of DIT 5 book Lecture - 15 : tips for Backprop

- We need it to work effectively
- Due to complex fn in neural net - Chance of getting stuck in local minima high



## ① Activation function

### Desirable Properties

Some neurons shouldn't explode

a) continuity of  $f(x)$ ,  $f'(x)$  <sup>to update weights we require  $f'(net)$</sup>

b)  $f(x)$  should saturate (max & min) <sup>we need this</sup>

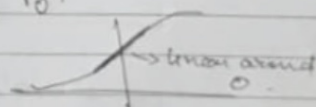
↳ to keep weights & activations in check  
↳ Also to interpret o/p of neuron as probability

c)  $f(x)$  monotonic

for  $J(\text{error})$  fn to be simple  $f$  must be simple

d)  $f(x)$  : linear around '0'

Sigmoid satisfy all these



## ② Parameters of Activation fn (Assume sigmoid is used)

$$f(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad \text{basic sigmoid}$$

We want some extra properties like the range of linearity & which part linear/non-linearity lies

$$f(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \Rightarrow \frac{1}{1 + e^{-b \cdot \text{net}}} \quad \text{b: control range of linearity}$$

ideally weights should be +ve & negative  
∴ o/p going from -ve to +ve is desirable

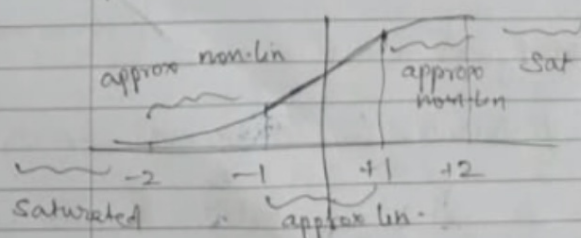
$$\therefore f(\text{net}) = \frac{1}{1 + e^{-b \cdot \text{net}}} - \frac{1}{2}$$

(Note: final o/p is wanted from 0 to 1)  
so maybe we want to do it for o/p

$$f(\text{net}) = 2a \left( \frac{1}{1 + e^{-b \cdot \text{net}}} - \frac{1}{2} \right)$$

multiplying this so that range changes from  $[-1/2, 1/2] \rightarrow [-a, a]$

$$f(\text{net}) = \frac{2a}{1 + e^{-b \cdot \text{net}}} - a$$



our goal ranges of lin & non-lin

$\therefore a = 1.716, b = 2/3$  gives desirable props

$$\max/\min (f'(net)) = +2/-2$$

$$\max/\min (f(net)) = +a/-a$$

$f(x) \Rightarrow$  also antisymmetric \*

$$f(-x) = -f(x)$$

- ③ Input values : if too large  $\Rightarrow$  saturate to  $-a/+a$   
 if too small  $\Rightarrow$  only non linearity becomes linear  
 $\therefore$  power of N:N lost - back to lin. perceptron

$\therefore$  shouldn't be too large / too small

\* Normalise : 0 mean, unit variance,  
 $\therefore \sim -2, +2$  range

- ④ Weights Initialization : should be in approp. range

$\sum_{j=1}^d w_j y_j$  } net should be in proper range

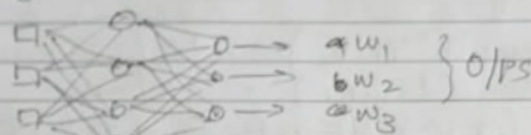
Note: don't set all weights to 0 (no learning happens, every eqn goes to 0)

So choose randomly

$\therefore W = U[-\tilde{w}, \tilde{w}]$  } uniform R.V

$$\star \tilde{w} = 1/\sqrt{d}$$

### ⑤ Target Value



lets assume  $w_1$  correspond to class 1  
 how is it encoded?  $w_1 \rightarrow$  High  $w_2, w_3 \rightarrow$  Low

$\Rightarrow w_1, w_2, w_3 = a, -a, -a$  . lets assume  
 in an attempt to make  $w_1$  go to  $a$   
 $\& w_2 \& w_3$  got  $-a$ . the net value at  
 $w_1$  goes to  $\infty$ , and at  $w_2, w_3$  go to  $-\infty$

$\therefore$  weights explode to push outputs to target values at saturation

\*  $\therefore$  ~~add~~ put target value less than saturation  
 eg:  $a = 1.716$  then put target at 1.5 or 1  
 both high & low should change  
 (ie high = 1.5 low = -1.5)

### ⑥ Training with Noise

- $\neq$  training samples small
- $\rightarrow$  add noise into data as training samples
- $\rightarrow$  helps regularising
- $\rightarrow$  if test data has noise, this helps
- $\rightarrow$  increases training set size,  $\rightarrow$  add to its vector

data augmentation

### ⑦ Manufacture new Data

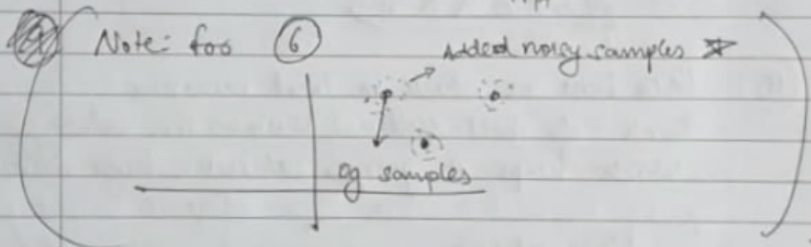
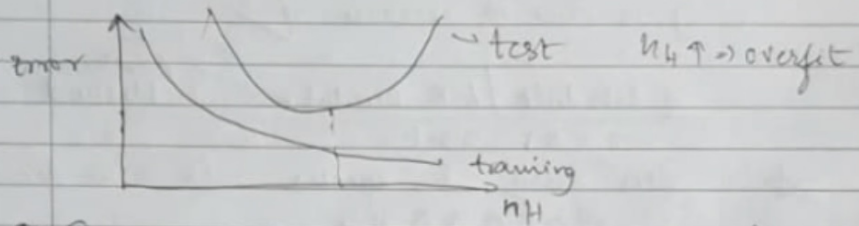
- $\rightarrow$  rotate, image process, etc on.
- $\rightarrow$  given one example get variants of  $\star$  sample
- $\rightarrow$  add as training data



⑧ # of hidden <sup>→ 104</sup> Units (Assume 3 layers)

i/p layer      hidden layer      o/p layer  
 $n_1$                        $n_H$                       # of classes

total 2 error method



⑨ Learning rate      too small  $\Rightarrow$  slow  
    too large  $\Rightarrow$  jump around. No learn

$\eta_{opt} \rightarrow$  optimal  $\eta$   
 e.g. (Adaptive Newton gradient descent method  
 can be used (approx for small regions &  
 take appropriate jump))

$\eta_{opt} = \left( \frac{\partial^2 J}{\partial w^2} \right)^{-1}$  } Computing is impractical \*  
    faster to do G.D.  
 $\eta_{max} = 2 \cdot \eta_{opt}$

1 step N.D.  $\approx$  10 step G.D. in time taken  
 also while 1 step N.D.  $\approx$  better than 1 step G.D.  
 10 step G.D. better than 1 step N.D.

Instead {  
 i. start with  $\eta = 0.1$   
 J decreases  $\Rightarrow$  reduce  $\eta$  \*  
 J is slow  $\Rightarrow$  increase  $\eta$

$\eta$  schedule / L-R schedule } Another Method -  
 $\eta = 0.1$  start  
 after every  $k$  epochs ( $k \rightarrow$  100 epochs)  
 $\eta = 0.9 \times \eta$

⑩ Don't get stuck in local minimum  
 think of ball rolling down  $\rightarrow$  has velocity  
 which keeps it going. (initial huge velocity)

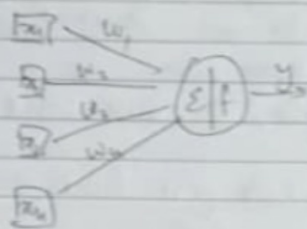
$\therefore$  Momentum

Oscillation possible  
 $\Delta w_{new} \Rightarrow (1-\alpha) \Delta w_{old} + \alpha \Delta w_{new}$   
 i.e. combination of prev. & new direction  
 (don't stop even if  $\Delta w_{new} = 0$ , continued in direction of  $\Delta w_{old}$ )

⑪ Weight Decay       $\lambda \approx 0.01, 0.05, 0.1$   
 $w \leftarrow w(1-\lambda)$  after every update  
 all numbers moved toward 0

→ let's say  $W_1$  is important ~~the~~ says pos.  
 G.D says push  $W_1$  up } G.D wins &  
 Weight decay says put it down } takes over

→ let's say  $W_2$  is not relevant  
 G.D don't push anywhere }  $W_2$  Decay wins  
 weight decay push to 0



∴ W. Decay makes  
all unused weights

☆ 〇

Useful weights determined by  
G.D.

★ Also Complexity of fn reduces as unwanted things go into linear region  $\Rightarrow$  Simplification

equivalent to  $J_{\text{new}} = J_{\text{original}} + \frac{2E}{n} \cdot w^t w$   
loss not for weight decay loss for

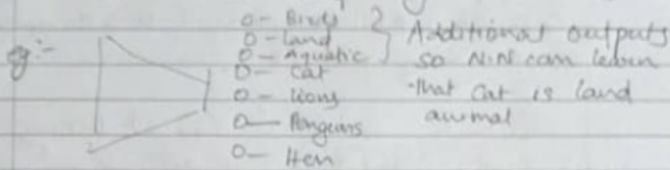
adding  $W^T W$  to  $J$  has to be minimized (known)  
 $\rightarrow$  we want to reduce values inside  $W$   
 i.e. size of  $W$  (minimizing  $W^T W$ )

2. this simple update eqn. to adding complex term to  $\mathcal{J}$  kept weights in check

Note:  $\frac{2\lambda}{\gamma}$  wt. w called regularization

12 Use of hints

Let's say we want to classify animals -  
add extra nodes corresponding to animals



when cat is higher, so will land output

★ Adding extra outputs to give hints to neural networks

We chop off extra ops at the end ~~\*~~

Unlike Normal: more than ~~one~~ one output high

(13) Loss fn / criterion fn

a) Square Error  $J(w) = \sum_k (t_k - z_k)^2$

5) Nankovska Sestra - 4

impulse response:  $J(\omega) = \sum_k (H_k - 2k) \delta$

c) Cross Entropy

c) Cross Entropy

$$J(w) = \sum_{k=1}^c t_k \ln(t_k / z_k)$$

⑭ # of hidden layers

(14) # of hidden layers  
# hidden layer  $> 1 \Rightarrow$  Deep Network  $\rightarrow$