# Parameter Calculation in DNNs. 19/02/25

batch size = 1



$1 \times 1 \times 1 \times 3072$

$1 \times \boxed{10} \times 1 \times 3072$

$1 \times 10 \times 1 \times 1$
——————
↓
output.

## Dimension of Convolution block :-

$$( \text{batch-size} \times \text{number of channels} \times \text{height} \times \text{width.} )$$



→ # of channels

h
w

5

5


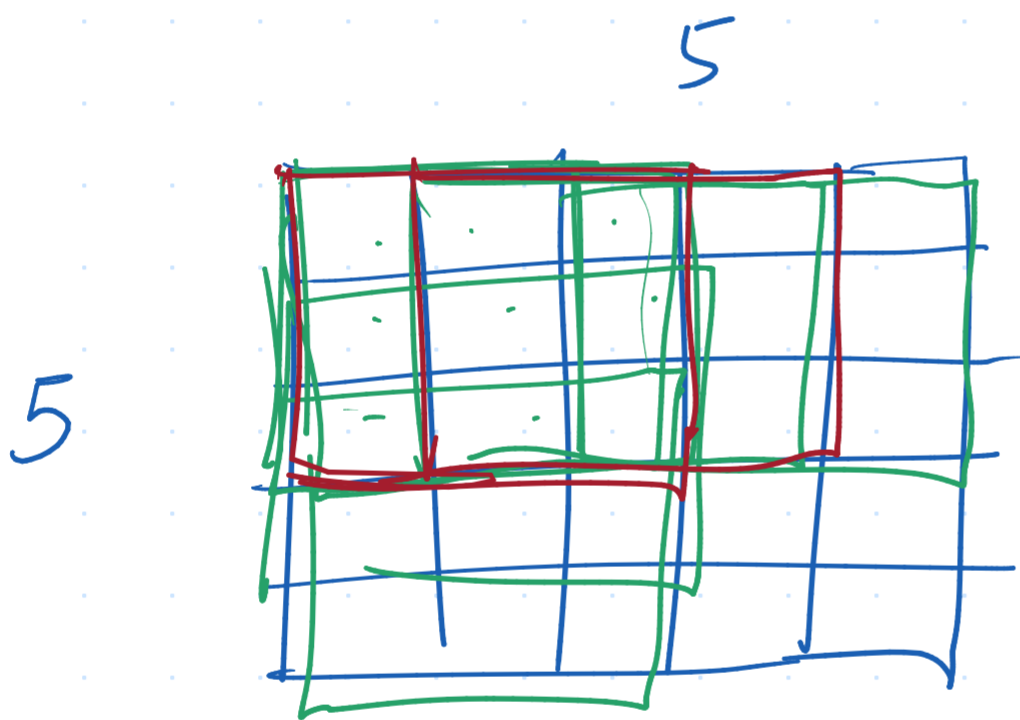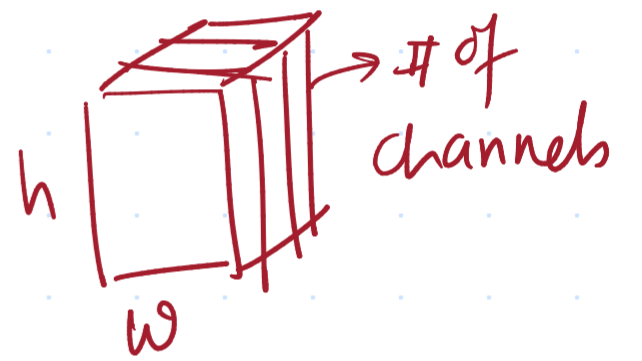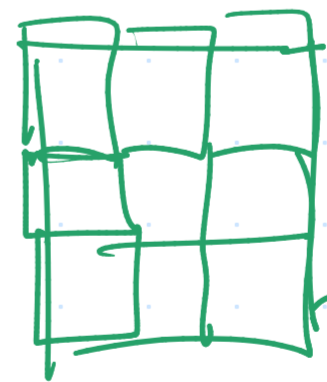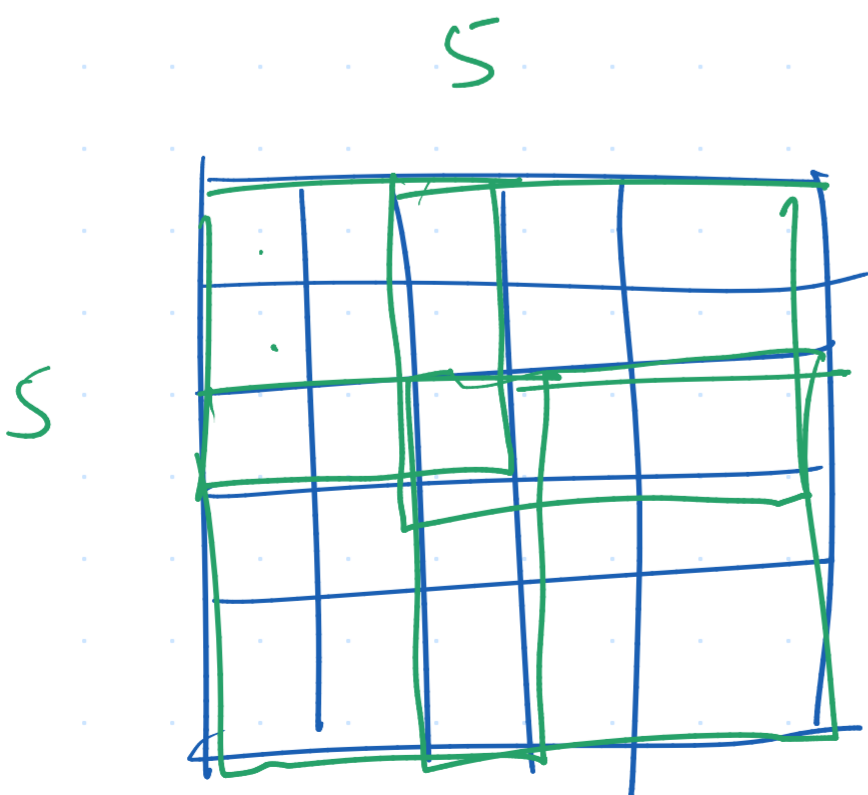
Image.

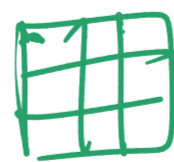1-channel → image.

$1 \times 1 \times H \times W$

3

3



convolution →

stride = 1

3×3.

5

S



3×3 kernel

with stride = 2

# Usage of Strides.

5000

4997

5000

3

3

stride 1 →

4998

4998

↓

feature map → increase with small stride.

↓ stride = 3.

$\left[\dfrac{5000}{3}\right]$

1667

1667

→ feature map size reduces.

CNN → mainly used to drastically reduce the size of the neural network.

3

3

kernel → # params ↗ 10

$3 \times 3 + 1$

↳ bias.

5

5

kernel → # params

$5 \times 5 + 1$ ↳ bias.

↙

26

$\dfrac{H \times W}{2}$

$O$
$O$  $O$
$O$  $O$
$O$
$O$

W

flatten  H

$1 \times 1 \times H \times W$

$\downarrow$ FCL / Dense Layer

# params $\dfrac{H \times W}{2} + H \times W$ ✓

flatten

## Receptive field.

$\xrightarrow{\quad 10 \quad}$

$\xrightarrow{\quad 10 \quad}$  # ⬭ 20 .

5

5

$\xrightarrow{\qquad\qquad}$  □  ^.

# params    26

$5 \times 5 + 1 \Rightarrow 26 .$

6

6

| 1 | 2 | 7 | 1 | 1 | 9 |
| 6 | 5 | 3 | 4 | 3 | 2 |
| 7 | 1 | 3 | 7 | 3 | 1 |

2×2 Maxpool.

| 6 | 7 | 9 |
| . | . | . |
| . | . | . |

feature Map.

$$\frac{6+2+5+1}{4}$$

(Average pooling)

3×3 Maxpool.

| 7 | 9 |
| × | × |

2×2  feature map

Dense

5

10.

10×5 + 10

Dense

Dense -

3.  3×10 + 3.

Sigmoid Activation

# parameter = ?

~~60~~
33

| 93 |

feature map.

$1 \times 3 \times 32 \times 32$

$\boxed{1 \times 3 \times 5 \times 5}$

3, 5

28
28

# params.

$3 \times 5 \times 5 + 1$.
$\Rightarrow 7571$
$\Rightarrow 76$

stride $= 1$

5

3

1+ 27

1

$32 - 5 \Rightarrow 27$.

convolve.

$2 \times 2 \times 2$

$2 \times 1 \times 1$.

4

1

2  C

| x1 | x2 |
|----|----|
| x3 | x4 |

convolution

| y1 | y2 |
|----|----|
| y3 | y4 |

$(x1 \, y1 + x2 \, y2 + x3 \, y3 + x4 \, y4.)$

# Convolution, correlation & filtering → Difference

O.S.

$$6 \times 3 \times 5 \times 5 \rightarrow \text{?} \qquad (\text{output ?})$$

ic    w

$$(\# \text{ parameters} = \text{?})$$

32

32

3

$bs \times \#C \times H \times W.$

$$(6 \times 3 \times 5 \times 5 + 6)$$

$$456$$

$$6 \times 1 \times 28 \times 28.$$

28

28

6.

$$6 \times 28 \times 28$$

# CNN Model Param Calculation.

$3 \times 3 \times 3 \times 3$

stride = 2
padding = 1

feature Map.

#param
$\Rightarrow 3 \times 3 \times 3 \times 3$
$+ 3$    #
$\Rightarrow 81 + 3 = \boxed{84}$

$3 \times 1 \times 4 \times 4.$

$3$
$2$
$27$
$\times 3$
$81$

padding = 3.
stride = 1

$\underset{O.S.}{3} \times \underset{I.C.}{3} \times \underset{h}{3} \times \underset{w}{3}$
kernel.

$1 + 7$
$0$
#
$84$

$1 \times 3 \times 8 \times 8$

$x - \underset{O}{M} \Rightarrow 3 \atop \Rightarrow 3.$

Maxpool.   brv

$3 + 3$

$1 \times 3 \times 8 \times 8$

$3 \times 4 \times 4.$
$48$

$1 \times 3 \times 4 \times 4$

$48 \times 10 + 10$

$48$

$O\ O\ O$ --- $O$

$\#$
$(10 + 10)$

$(48 \times 10 + 10)$

$490$

$10$ $O$ $O\ O\ O$ $\boxed{10}$

$O$

$1 \rightarrow$ $(10 \times 1 + 1)$

$11$

---

Backprop & Gradient Descent

Greedy Algorithm

It always moves in the direction of reducing errors.

$O$
$O$
$O$
$O$

In N.N. one of the local optima is as good as the global optima

initialization will matter if you are using Gradient descent

Adam → $\vdots$ → One of the local optima → which is as good as global optima.

Weight update rule to minimize errors in loss functions.

$$\Delta w_{ji} = -\eta \frac{\partial E}{\partial w_{ji}}$$



98·3%0·01

98·37%0·000%

98·43%0·0093

$\eta$ → Learning rate

$E$ → Loss function

$w_{ji}$ → weight of connection from the $i$th neuron to the $j$th neuron.

$\{98\cdot3, \ 98\cdot37, \ 98\cdot43, ----\}$

$\{98\cdot36 \pm 0\cdot07\}$

report

Always run for atleast 10-15 times & report the mean & standard deviation on the fist dataset

Linear Regression

Convex → Always converge to global minima

$E\uparrow$

whatever may be the initialization



$E\uparrow$

$w_{ji}$

$\frac{\partial E}{\partial w_{ji}} < 0$ , $\Delta w_{ji} \Rightarrow +ve$.

A

$\Delta w_{ji} = -ve$.

$\frac{\partial E}{\partial w_{ji}} > 0$,

$\frac{\partial E}{\partial w_{ji}} = 0$.

Stationary.

$w_{ji}$

$x$-axis

We will now derive the single sigmoid Neuron's cross entropy loss from a single data point

$$\frac{\partial E}{\partial w_1} = \left(\frac{\partial E}{\partial \theta}\right) \cdot \left(\frac{\partial \theta}{\partial net}\right) \cdot \left(\frac{\partial net}{\partial w_i}\right)$$

$$net = \sum_i x_i w_i$$



$$E = -t \log \theta - (1-t) \log(1-\theta)$$

cross entropy → binary

$$\frac{\partial E}{\partial \theta} = -\frac{t}{\theta} - \frac{(1-t)}{(1-\theta)}(-1)$$

[ Small $\theta$ ]

$$= \frac{-t(1-\theta) + (1-t)\theta}{\theta(1-\theta)} = \frac{-t + t\theta + \theta - t\theta}{\theta(1-\theta)}$$

$$= \frac{-t + \theta}{\theta(1-\theta)}$$

$$\boxed{\frac{\partial E}{\partial \theta} = -\frac{(t-\theta)}{\theta(1-\theta)}}$$

$$\boxed{\frac{\partial \theta}{\partial net} = \theta(1-\theta)}$$

$$\boxed{\frac{\partial net}{\partial w_i} = x_1}$$

$$\frac{\partial E}{\partial w_1} = -\frac{(t-\theta)}{\theta(1-\theta)} \, \theta(1-\theta) \cdot x_1 = -(t-\theta)x_1$$

$$\left[ \; -(t-\theta)x_m, \; -(t-\theta)x_{m-1}, \; \cdots \; -(t-\theta)x_1, \; (t-\theta)x_0 \right.$$

$$\boxed{\Delta w_1 = -\eta \frac{\partial E}{\partial w_1} = \eta(t-0)x_1}$$

$$\boxed{\Delta w_i = -\eta \boxed{\frac{\partial E}{\partial w_i}} = \eta(t-0)x_i} \implies \underline{\text{Derivative}}$$

Multiple Neurons $\rightarrow$ m output layer $\rightarrow$ Softmax & Cross Entropy

loss $\rightarrow$ illustrated with 2 neurons.



$$O = \langle o_1, o_0 \rangle$$

$$NET = \langle net_1, net_0 \rangle$$

$$Softmax_i = \frac{e^{net_i}}{\sum_j e^{net_j}}$$

$$O_1 = \frac{e^{net_1}}{e^{net_0} + e^{net_1}} \qquad O_0 = \frac{e^{net_0}}{e^{net_0} + e^{net_1}}$$

$$\frac{\partial O}{\partial Net} = \begin{bmatrix} \frac{\partial O_0}{\partial net_0} & \frac{\partial O_1}{\partial net_0} \\ \frac{\partial O_0}{\partial net_1} & \frac{\partial O_0}{\partial net_1} \end{bmatrix} = \begin{bmatrix} O_0(1-O_0) & -O_0 O_1 \\ -O_1 O_0 & O_1(1-O_1) \end{bmatrix}$$

jacobian

Refer to the derivative of softmax function.

for same $\rightarrow O_i(1-O_i)$

for diff $\rightarrow -O_i O_j$

Jacobian $\rightarrow$ first order partial derivative.

Hessian $\rightarrow$ 2nd order partial derivative.

m-tensor $\rightarrow$ Higher order partial derivative.

3-tensor. , 4-tensor. . . .

$\rightarrow$ Binary cross entropy

$$E = -t_1 \log O_1 - (1-t_1) \log(1-O_1)$$

$$\frac{\partial E}{\partial w_{11}} = -\frac{t_1}{O_1} \frac{\partial O_1}{\partial w_{11}} - \frac{t_0}{O_0} \frac{\partial O_0}{\partial w_{11}}$$

$$\frac{\partial O_1}{\partial w_{11}} = \frac{\partial O_1}{\partial net_1} \cdot \frac{\partial net_1}{\partial w_{11}} + \frac{\partial O_1}{\partial net_0} \cdot \frac{\partial net_0}{\partial w_{11}} \nearrow 0.$$

$$= O_1(1-O_1) \cdot x_1 + -O_1 O_0 \cdot 0$$

$$= O_1(1-O_1) \cdot x_1.$$

$$\frac{\partial o_0}{\partial w_{11}} = \frac{\partial o_0}{\partial net_1} \cdot \frac{\partial net_1}{\partial w_{11}} + \frac{\partial o_0}{\partial net_0} \cdot \frac{\partial net_0}{\partial w_{11}}^{\nearrow 0}$$

$$= - o_0 o_1 \, x_1 \quad + \quad - o_0(1-o_0) \times 0^{\nearrow 0}$$

$$\frac{\partial E}{\partial w_{11}} = - \frac{t_1}{o_1} \cdot \frac{\partial o_1}{\partial w_{11}} - \frac{t_0}{o_0} \cdot \frac{\partial o_0}{\partial w_{11}}$$

$$= - \frac{t_1}{o_1} \cdot o_1(1-o_1) \cdot x_1 - \frac{t_0}{o_0} (- o_0 o_1) \cdot x_1$$

$$= - t_1(1-o_1) \cdot x_1 + t_0 o_1 x_1$$

$$= x_1(-t_1 + o_1) = -(t_1 - o_1) \cdot x_1$$

$$\boxed{\Delta w_{11} = -\eta \, \frac{\partial E}{\partial w_{11}} = \eta \, (t_1 - o_1) x_1}$$