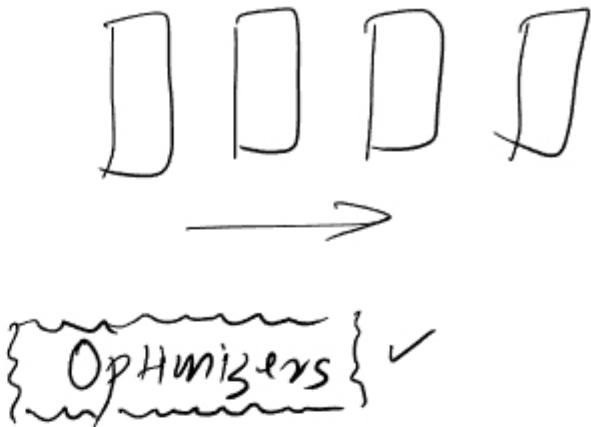


Introduction

05 July 2022 09:57

Performance of NN

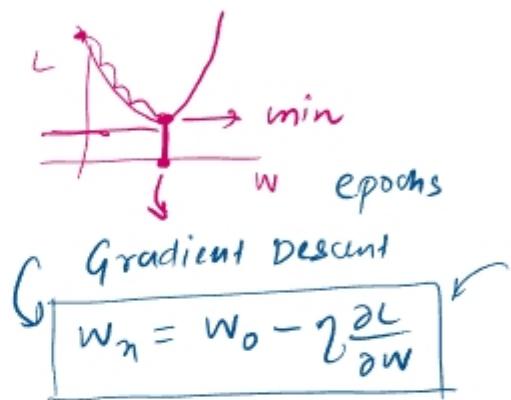
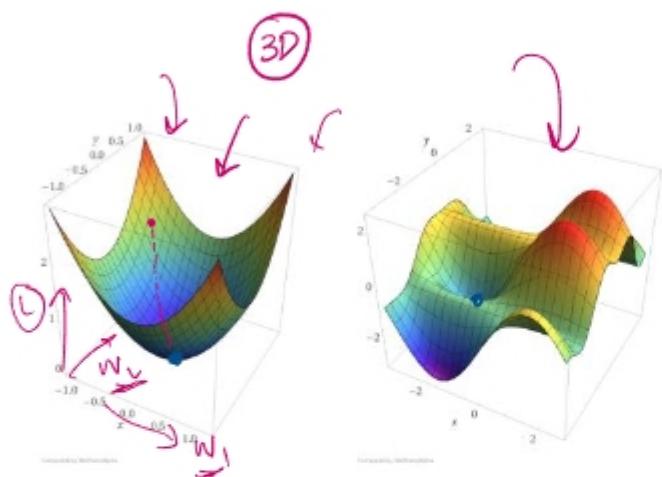
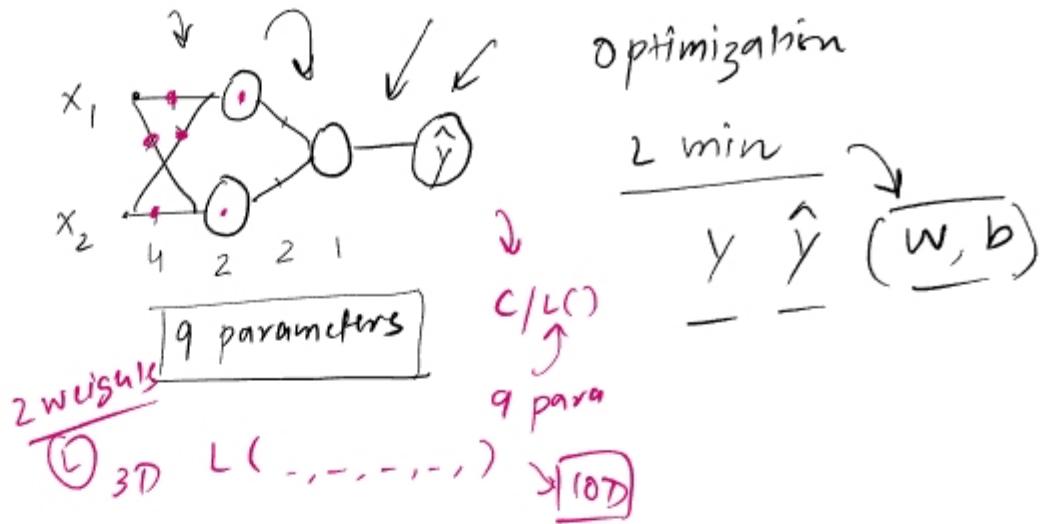
↓ how the speed the
training



- Weigh init
- Batch Norm
- Activation function

Role of Optimizer

05 July 2022 10:01



Types of Optimizers

05 July 2022 10:02

- Batch GD
- Stochastic GD
- Mini batch GD

10

epochs = 10

$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

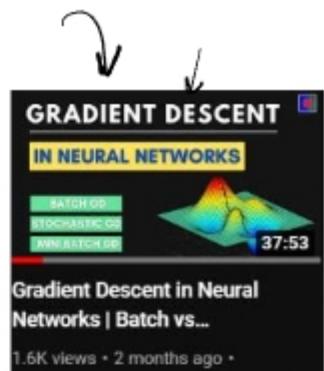
500 rows → pred

↓
loss → weight update

$[10 \times 500] \rightarrow 5000$

batch size = 100

5 batch → 10×5 times

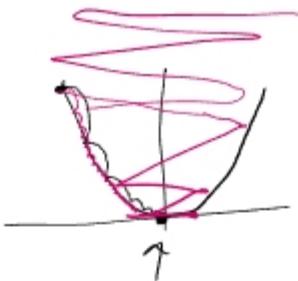


Challenges

05 July 2022 10:02

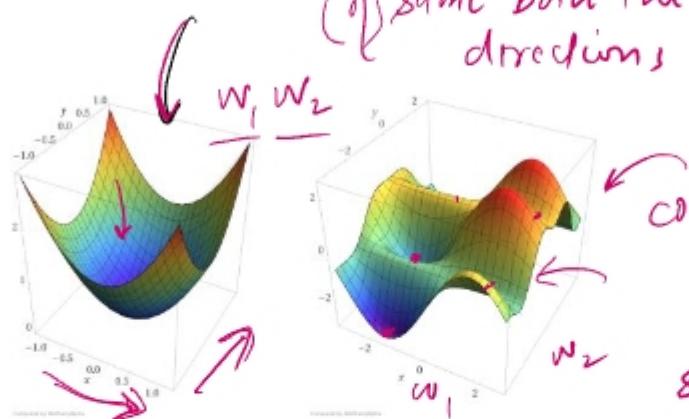
1) learning rate

$$w_n = w_0 - \eta \frac{\partial L}{\partial w_0}$$



2) learning rate scheduling → pre define

3)



⑨ weight's and bias' ⑨ sep learn

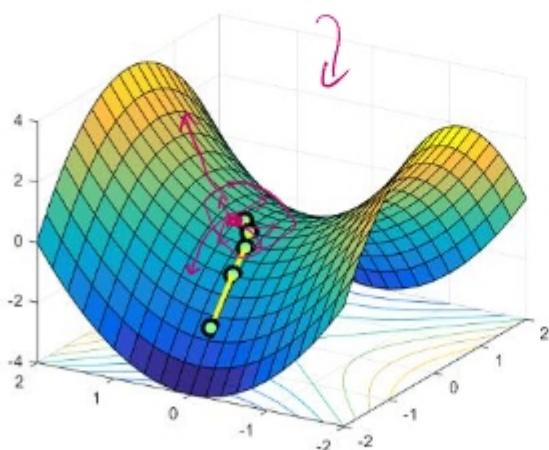
10-D → 10 dimensions

complex minima

mm w^T
sub optimum

4) local minima

5) saddle point



$$\frac{\partial L}{\partial w} = 0$$

$$w_n = w_0 - \eta \left[\frac{\partial L}{\partial w} \right]$$

What next?

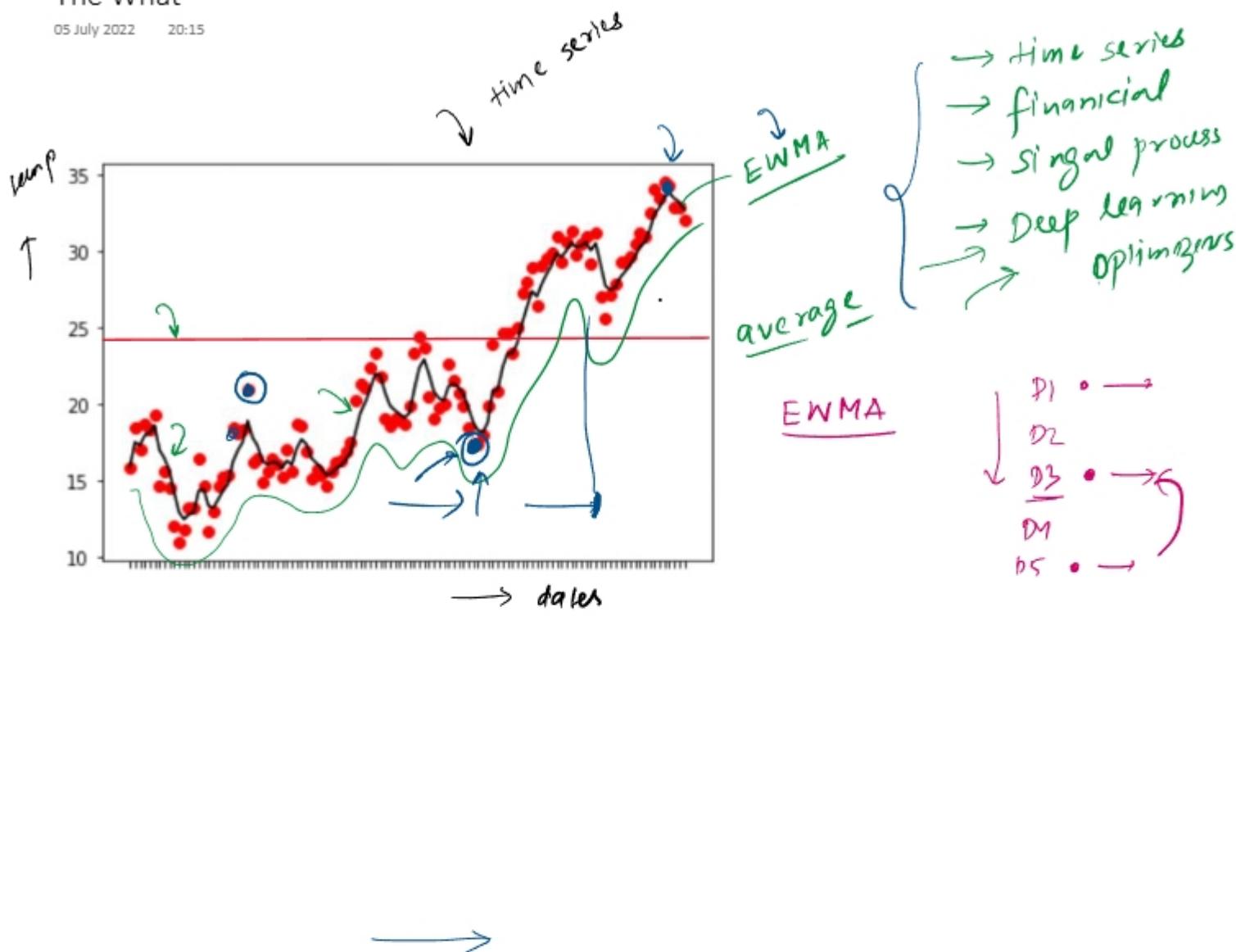
05 July 2022 10:02

- 1) Momentum
- 2) Adagrad
- 3) NAG
- 4) RMSprop
- 5) Adam

Exponentially
Weighted
Moving
Average

The What

05 July 2022 20:15



Mathematical Formulation

05 July 2022 20:16

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$\beta \rightarrow 0 < \beta < 1$

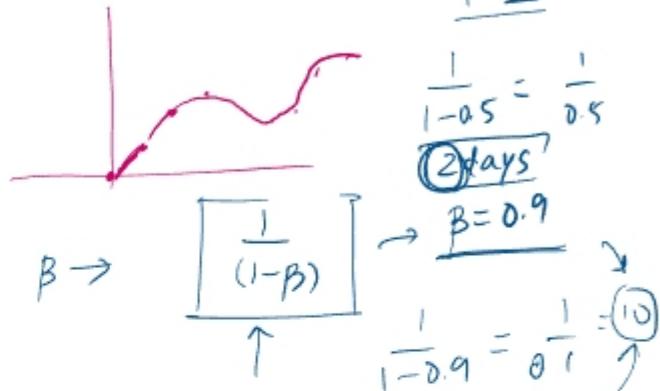
$$V_0 = 0 \quad \theta_0$$

$$\begin{aligned} V_1 &= 0.9 \times V_0 + 0.1 \times 13 \\ &= 1.3 \end{aligned}$$

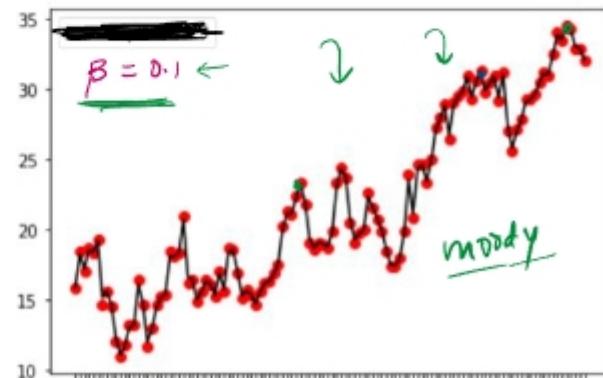
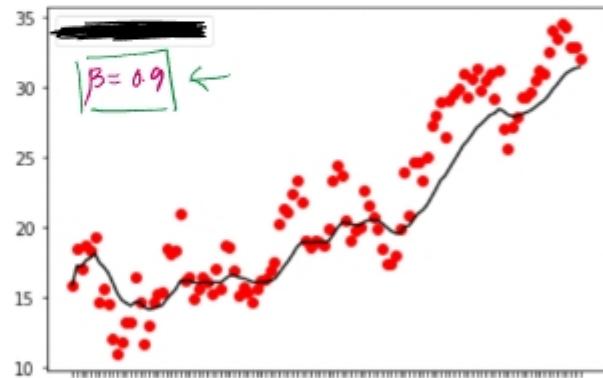
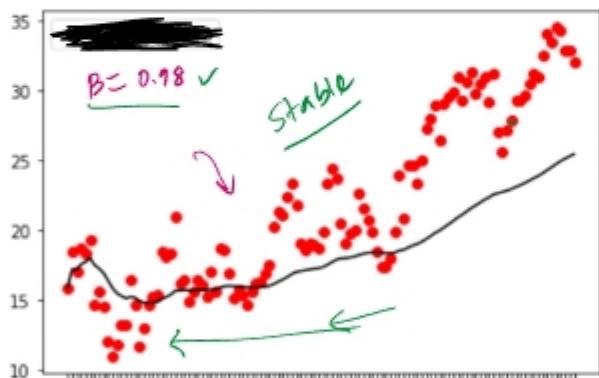
$$V_2 = 0.9 \times 1.3 + 0.1 \times 17 =$$

Index	temp (θ)
D1	25
D2	13
D3	17
D4	31
D5	43

$$\beta = 0.9$$



Effect of β



Mathematical Intuition

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

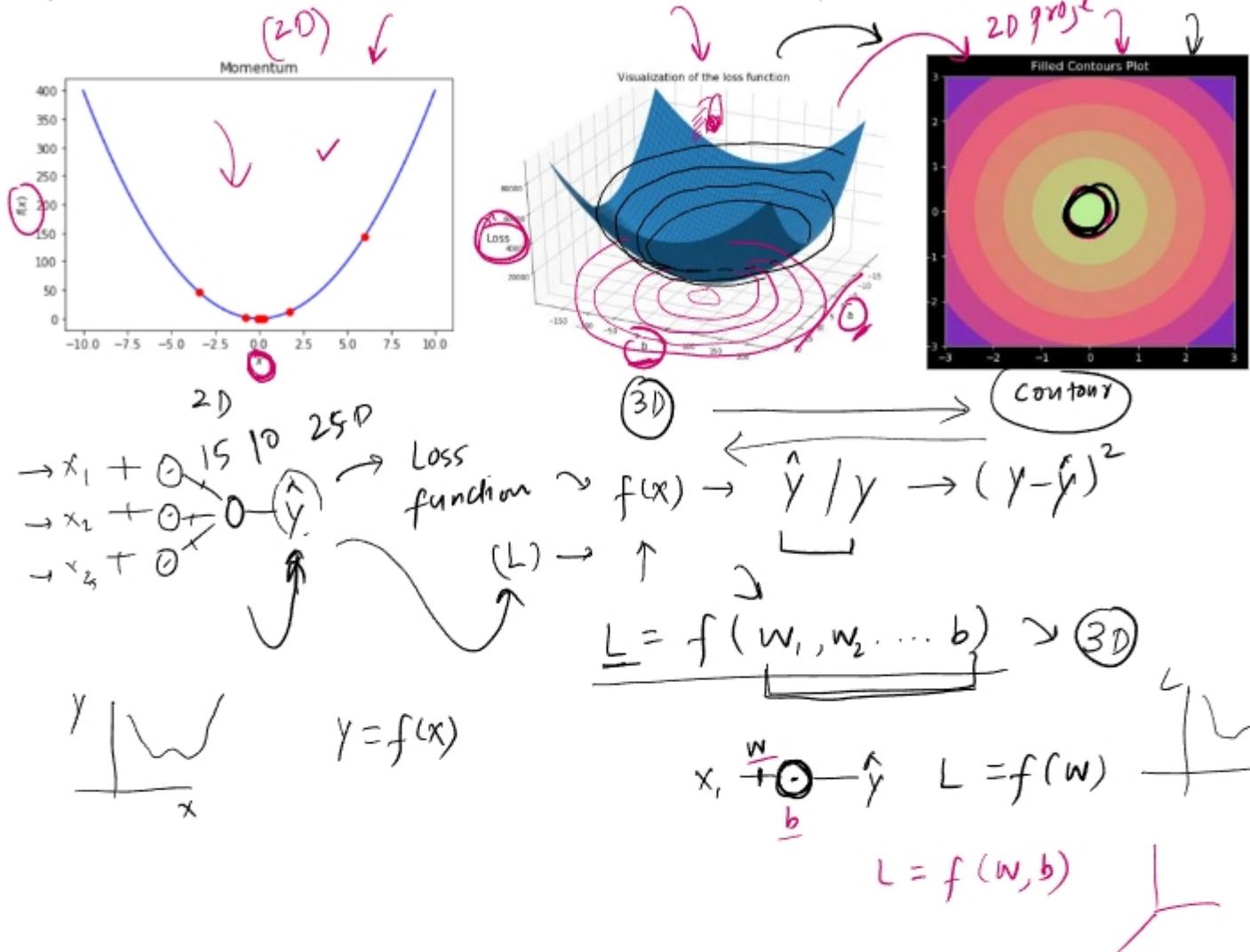
$$\left. \begin{aligned} V_4 &= \beta V_3 + (1-\beta) \theta_4 \\ V_4 &= \beta^3 (1-\beta) \theta_1 + \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_3 + (1-\beta) \theta_4 \end{aligned} \right\}$$

$$\left\{
 \begin{array}{l}
 V_t = \beta V_{t-1} + (1-\beta) \theta_t \\
 V_0 = 0 \\
 V_1 = (1-\beta) \theta_1 \\
 V_2 = \beta V_1 + (1-\beta) \theta_2 \\
 = \frac{\beta (1-\beta) \theta_1 + (1-\beta) \theta_2}{\beta V_2 + (1-\beta) \theta_2} \\
 V_3 = \frac{\beta^2 (1-\beta) \theta_1 + \beta (1-\beta) \theta_2 + (1-\beta) \theta_3}{\beta^3 (1-\beta) \theta_1 + \beta^2 \theta_2 + \beta \theta_3 + \theta_4}
 \end{array}
 \right.$$

$\beta^3 < \beta^2 < \beta$ $0 < \beta < 1$

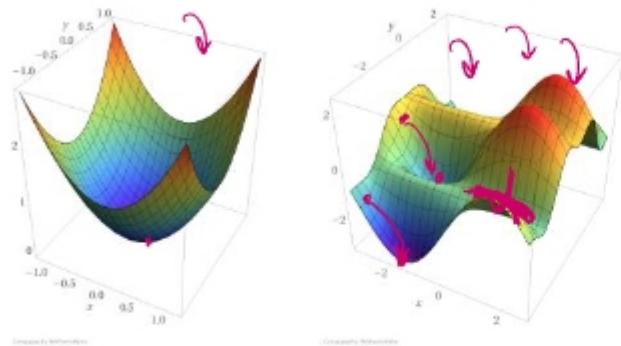
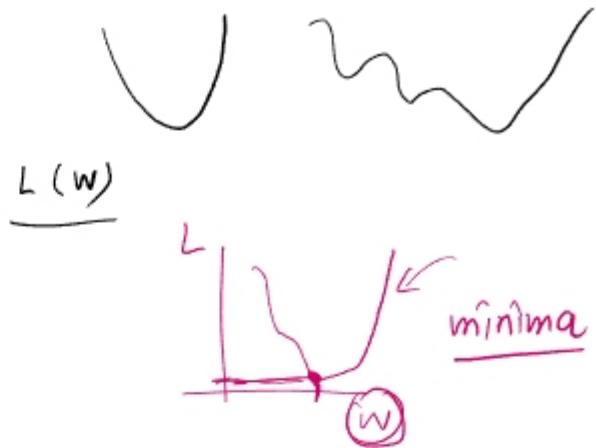
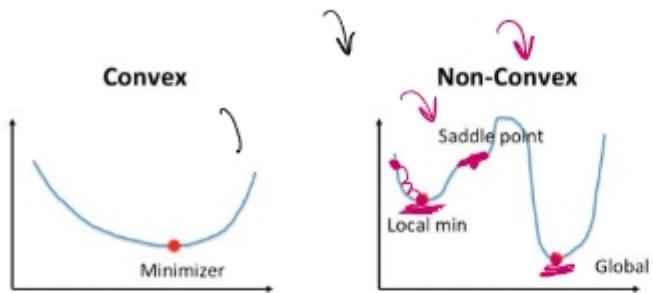

Understanding Graphs

20 July 2022 12:03

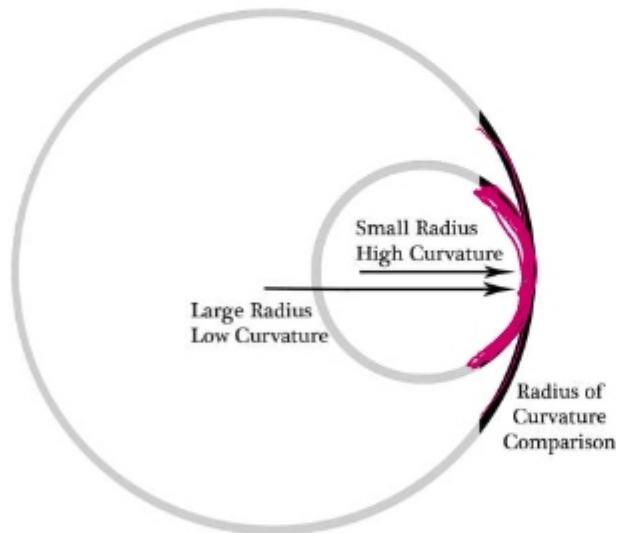


Convex Vs Non-Convex Optimization

20 July 2022 13:06

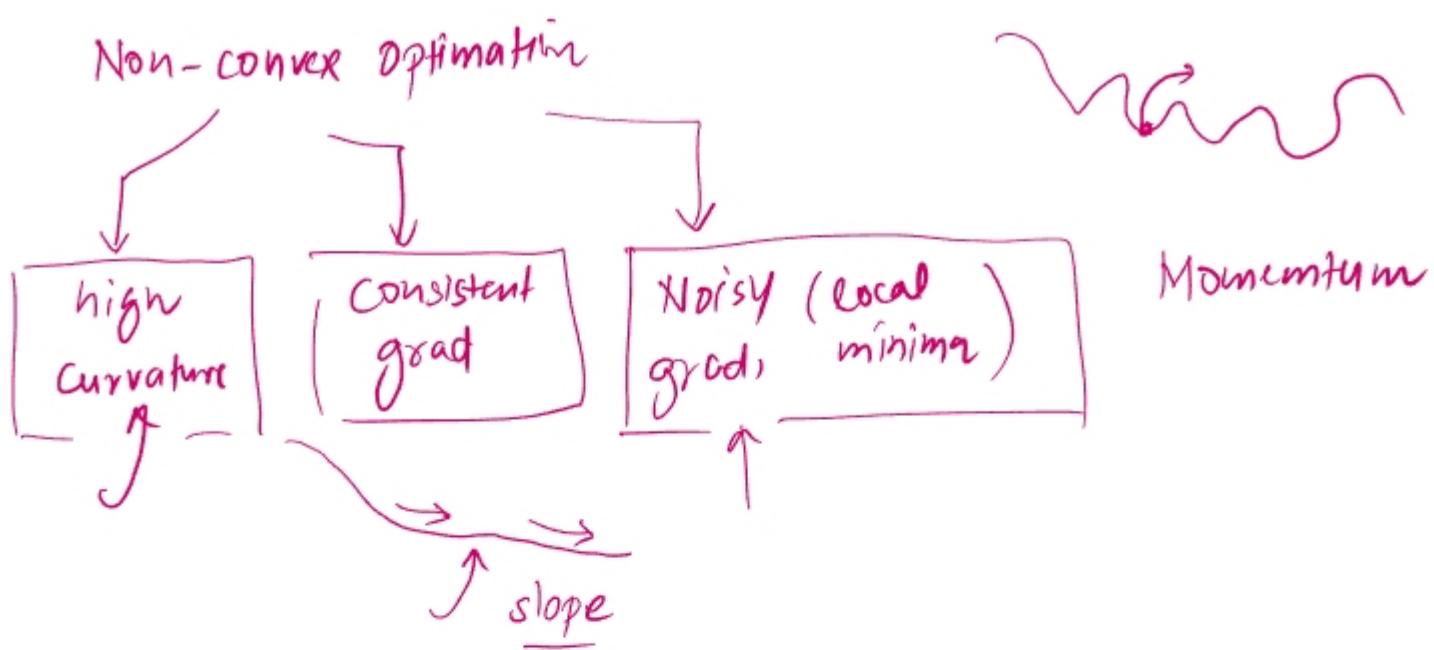


- 1) local minimum
- 2) saddle point
- 3) High curvature



Momentum Optimization - The Why?

20 July 2022 14:28

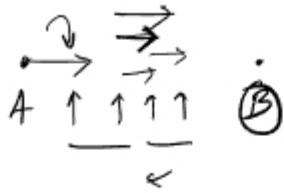


Momentum Optimization - The What?

20 July 2022 14:33

Speed

Common



prev grad \rightarrow

$$\begin{matrix} \rightarrow \\ \downarrow \\ mv \end{matrix}$$

$m \rightarrow$ unit mass
velocity \rightarrow previous history update



Momentum Optimization - Mathematics(The How?)

20 July 2022 14:56

$$\underline{w_{t+1}} = \underline{\dot{w}_t} - \boxed{\eta \nabla w_t}$$

}

$v \rightarrow \text{velocity}$

$w_{t+1} = w_t - v_t \rightarrow$

acceleration

momentum

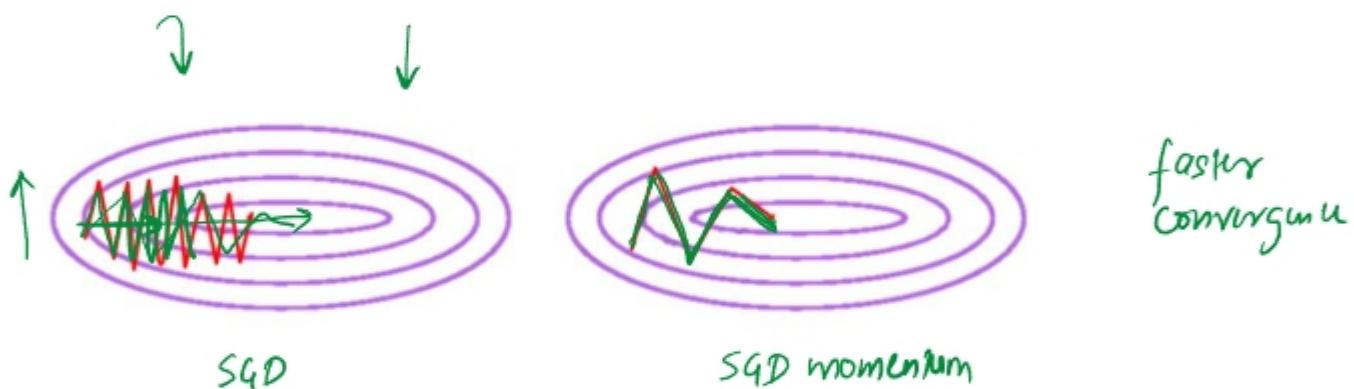
history of past velo- *use*

$v_t = \boxed{\beta * v_{t-1} + \eta \nabla w_t}$

$0 < \beta < 1$

v_{t-2}

v_{t-3}



Effect of beta

20 July 2022 15:09

$$w_{t+1} = w_t - v_t$$

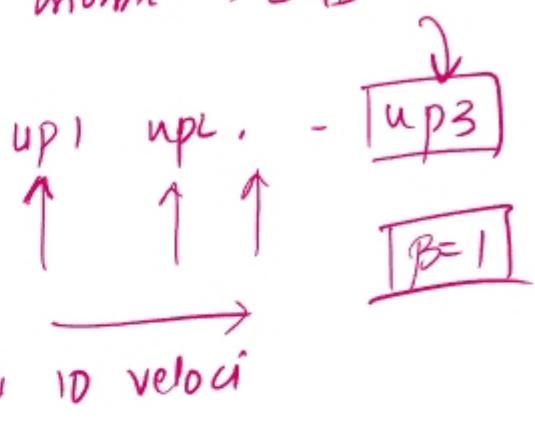
$$v_t = \beta v_{t-1} + \eta \nabla w_t$$

$\boxed{\beta=0}$

SGD

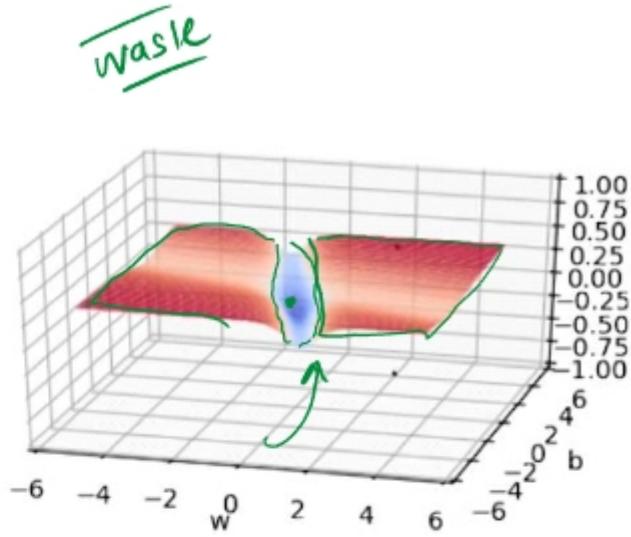
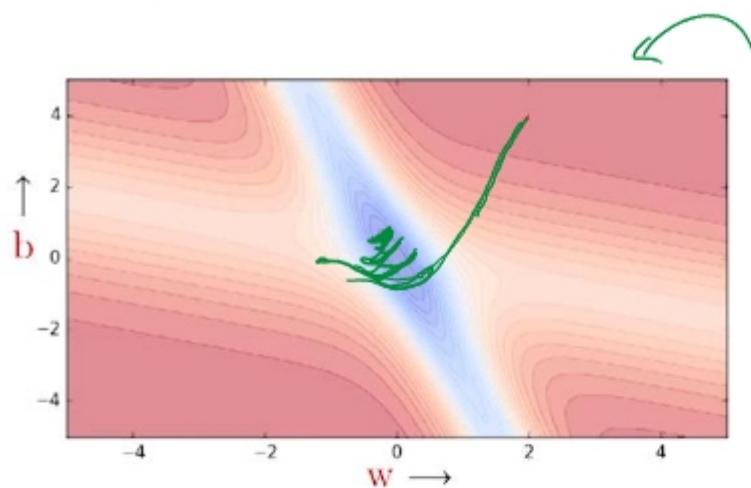
$\beta=0$ momen \rightarrow SGD

$$\frac{1}{1-\beta} = \frac{1-10}{0.1} \frac{1}{1-\beta} \quad \begin{matrix} \text{decaying factor} \\ 0.9 \\ \text{current} \\ \text{avg past} \end{matrix}$$



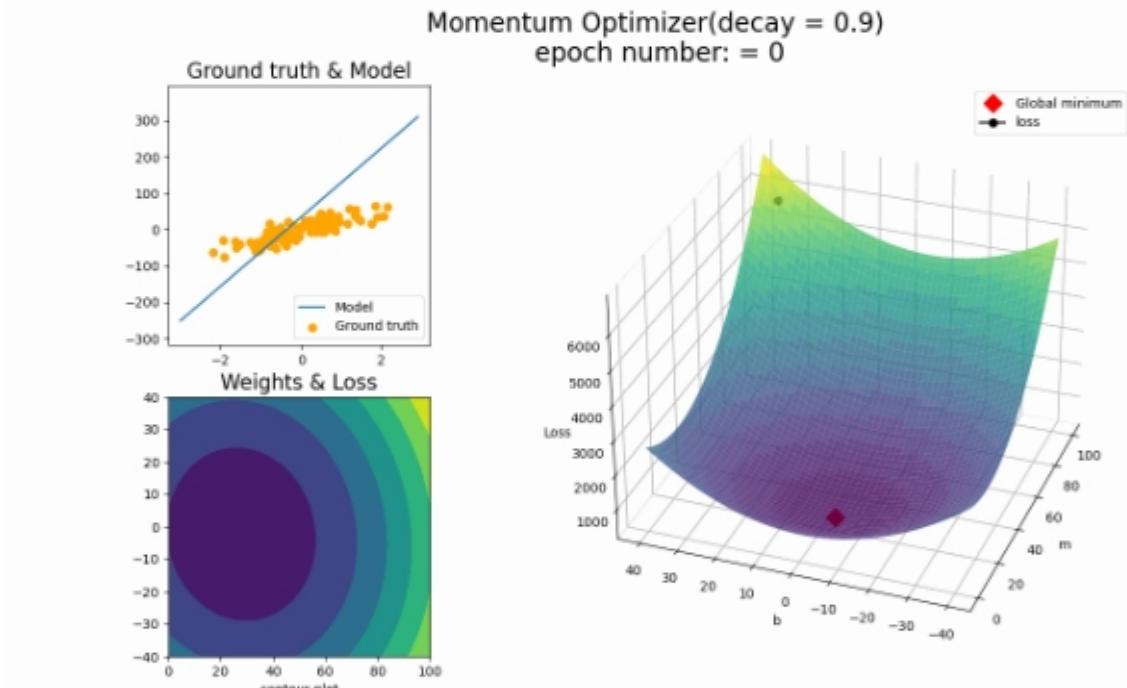
Problem with momentum optimization

20 July 2022 15:16



waste
decay 0.9

Momentum Optimizer(decay = 0.9)
epoch number: = 0



Demo

24 July 2022 13:17

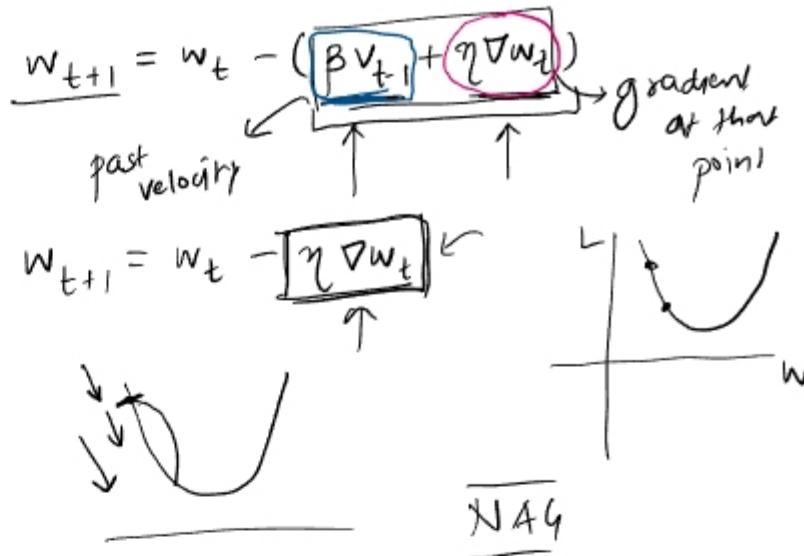
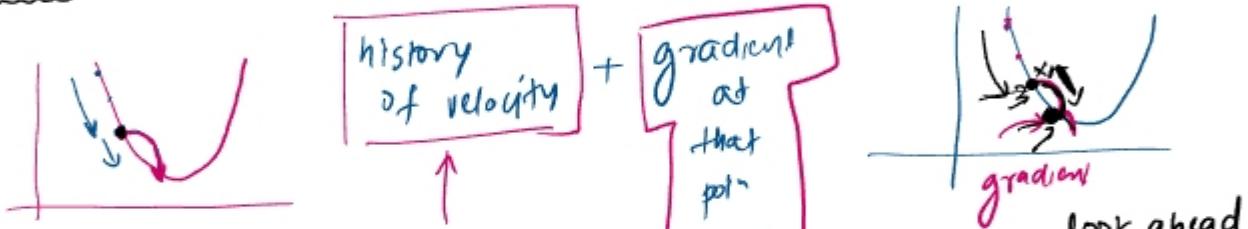
Momentum

$$w_{t+1} = w_t - v_t$$

where

$$v_t = \beta v_{t-1} + \eta \nabla w_t$$

$\beta \rightarrow$ decay factor

Nesterov Accelerated Gradient *better*

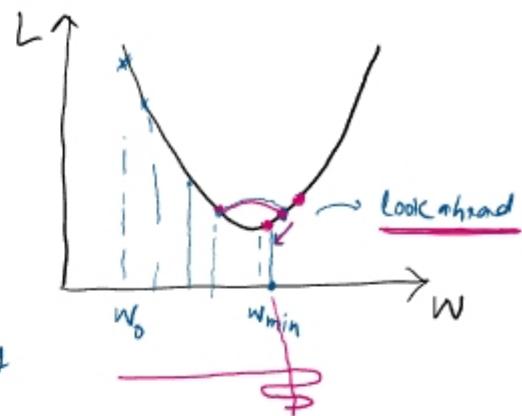
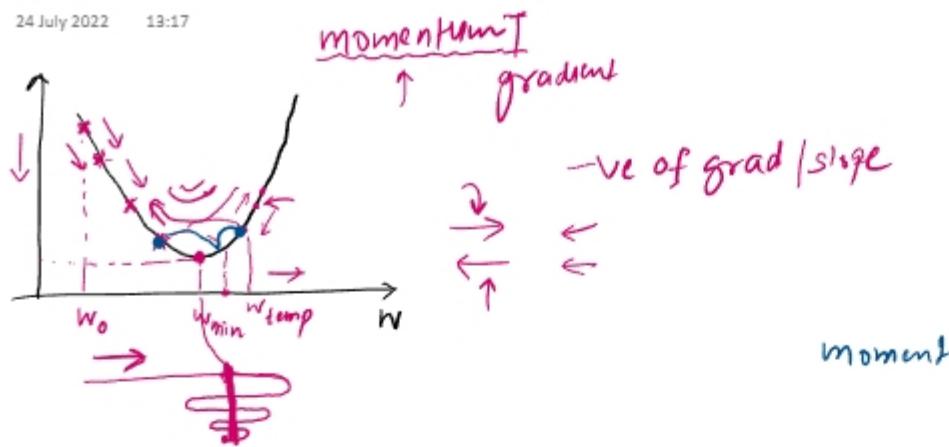
$$\left\{ \begin{array}{l} w_{\text{la}} = w_t - \beta v_{t-1} \\ v_t = \beta v_{t-1} + \eta \nabla w_{\text{la}} \\ w_{t+1} = w_t - v_t \end{array} \right.$$

$v_t = \frac{(w_t - w_{\text{la}}) + \eta \nabla w_{\text{la}}}{\beta}$

look ahead

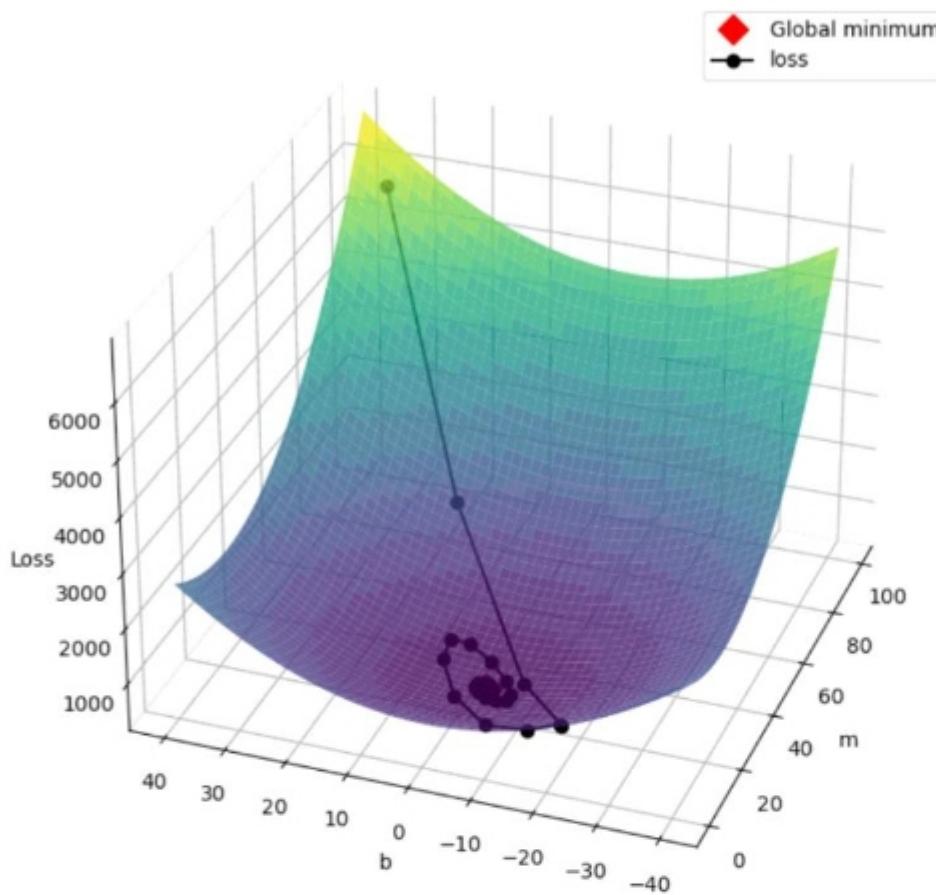
Geometric Intuition

24 July 2022 13:17



Disadvantage

24 July 2022 13:18



Disadvantage

NAG \rightarrow oscillation
 \downarrow
dampen
local minima

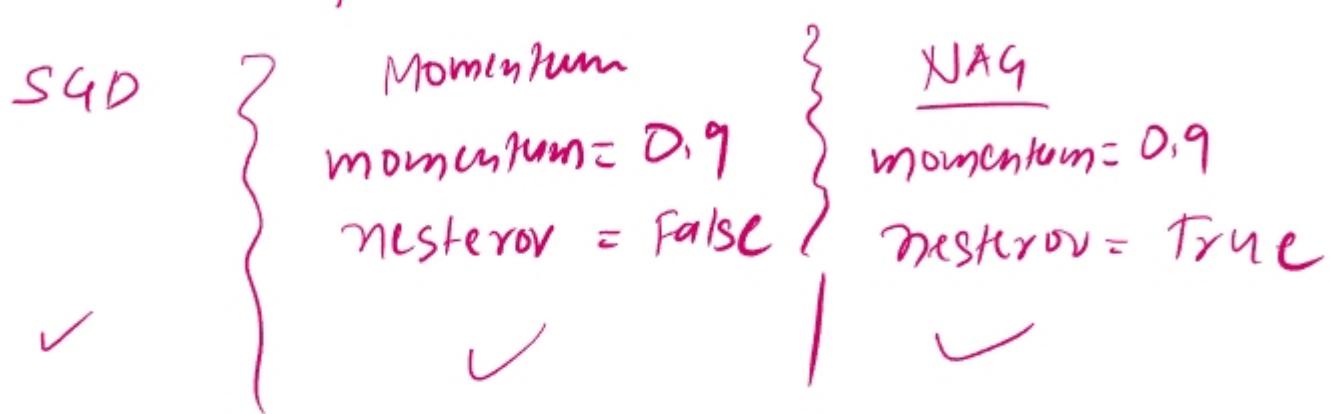


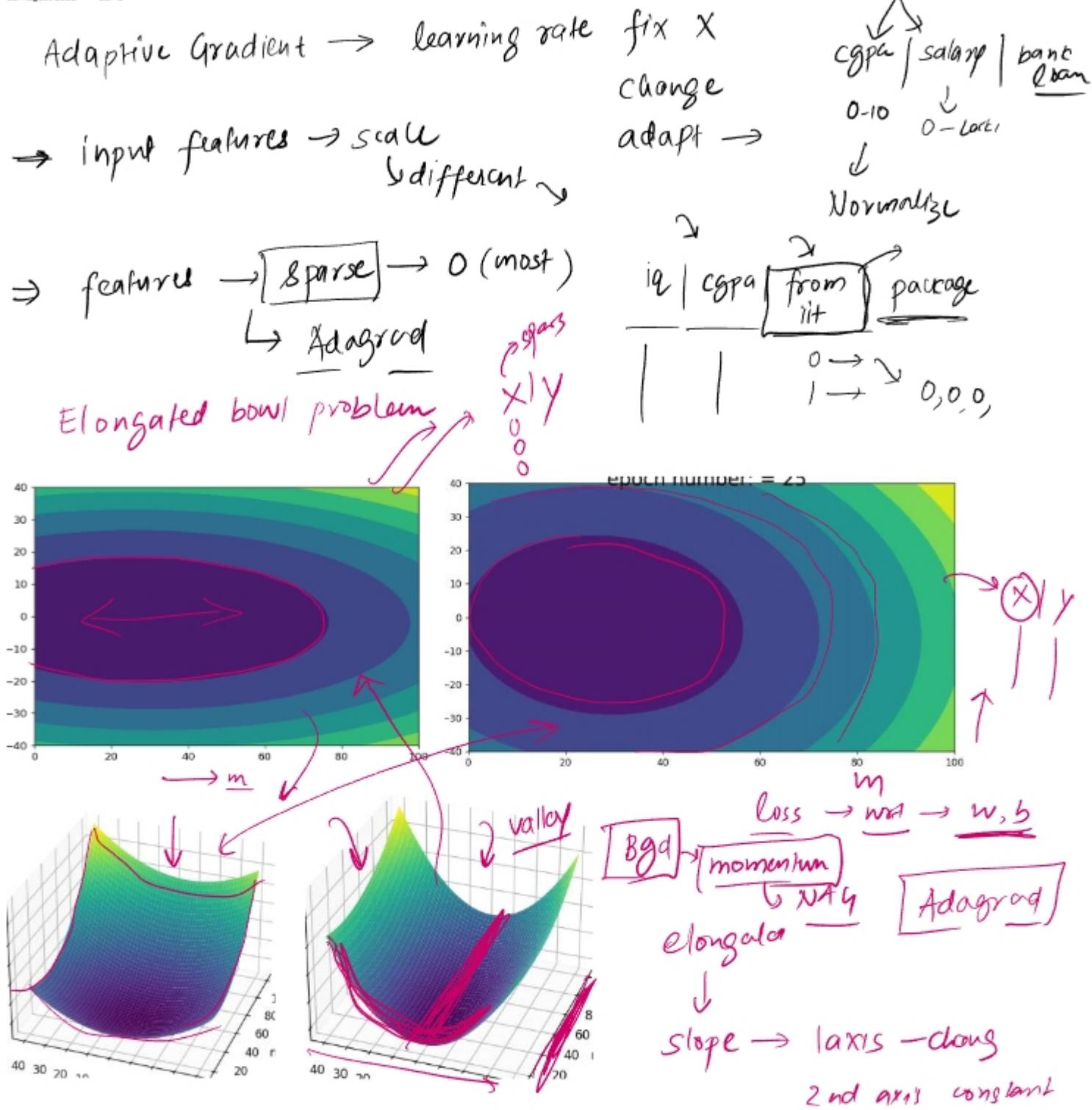
Keras Code

24 July 2022 13:18

2

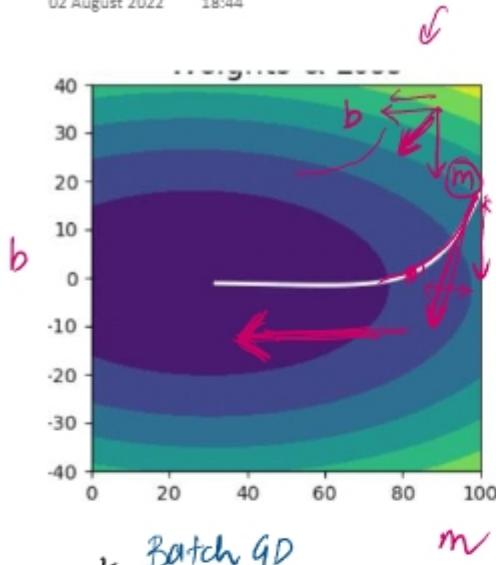
```
tf.keras.optimizers.SGD(  
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs  
)
```



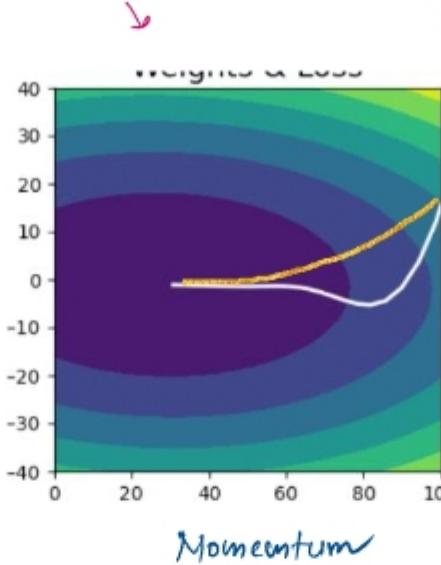


How optimizers behave (Why?)

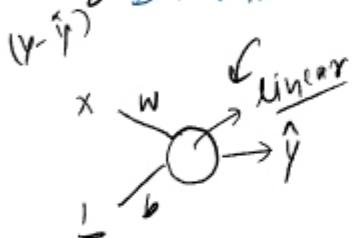
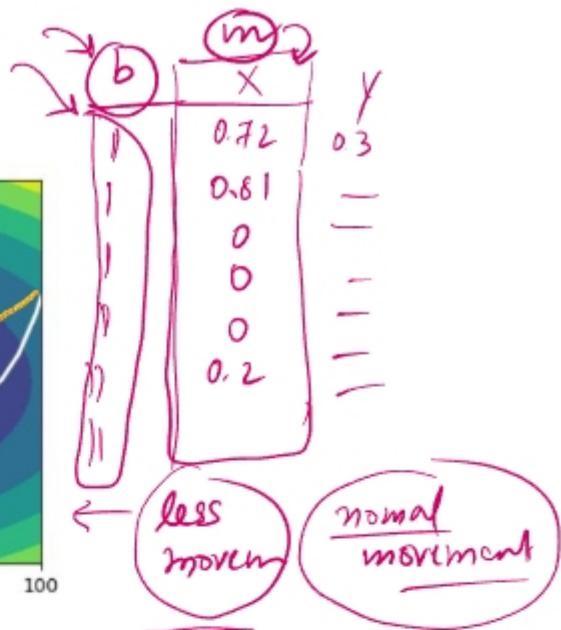
02 August 2022 18:44



Batch GD



Momentum



BGD → sparse

for i in epochs:

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial w} \\ = -2(y - \hat{y}) \times x \quad \text{add}$$

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$b = b - \eta \frac{\partial L}{\partial b}$$

→ big update in every epoch

$$\frac{\partial L}{\partial b} = [-2(y - \hat{y})] \times 1 \quad \text{big num}$$

sparse col →
small non sparse
↓ normal/big

Adagrad mathematics + Intuition

02 August 2022 18:44

$w = w - \eta \frac{\partial L}{\partial w}$ update

Adagrad

$b = b - \eta \frac{\partial L}{\partial b}$ optimization
lr → same

Adagrad → diff learning rates
big → lr → small
small → lr → big

$v_t = v_{t-1} + (\nabla w_t)^2$ past gradients sum (squared)
epoch

$\frac{\partial L}{\partial w}$ v_t $(\frac{\partial L}{\partial w})^2$ $(\frac{\partial L}{\partial b})^2$

$w_{t+1} = w_t - \eta \nabla w_t$ η small $w_{t+1} = w_t - \eta \Delta w_t$

$v_t + \epsilon$ η $\epsilon = \text{small num}$ $(-\text{ve})$ $(+\text{ve})$

$v_t = 0$

sparse data

Adagrad Demo

02 August 2022 18:44

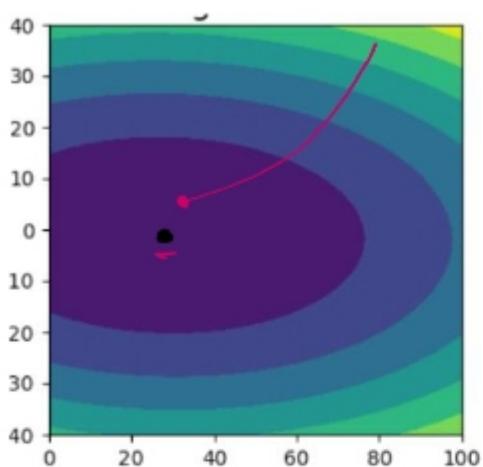
Keras Implementation

02 August 2022 18:44

Disadvantage

02 August 2022 18:43

NN \rightarrow use $\times \curvearrowright$ linear regress



$$\frac{\eta}{\sqrt{v_t}} v_t$$

$\xrightarrow{\text{small}} \text{decrease}$
 $\xrightarrow{\text{past updates}} \text{gradient}$
 $\xrightarrow{\text{epochs}}$
 $\xrightarrow{\text{update small}} \text{global minima}$
 $\xrightarrow{\text{RMSprop}}$
 $\xrightarrow{\text{Adam}}$

The Why?

03 August 2022 14:03

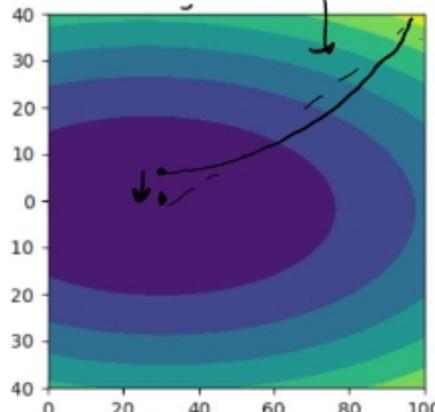
RMSprop → Root Mean square Prop

↑ improvement

→ Adagrad ← sparse data

BGD → momentum

Disadvantage → learning rate

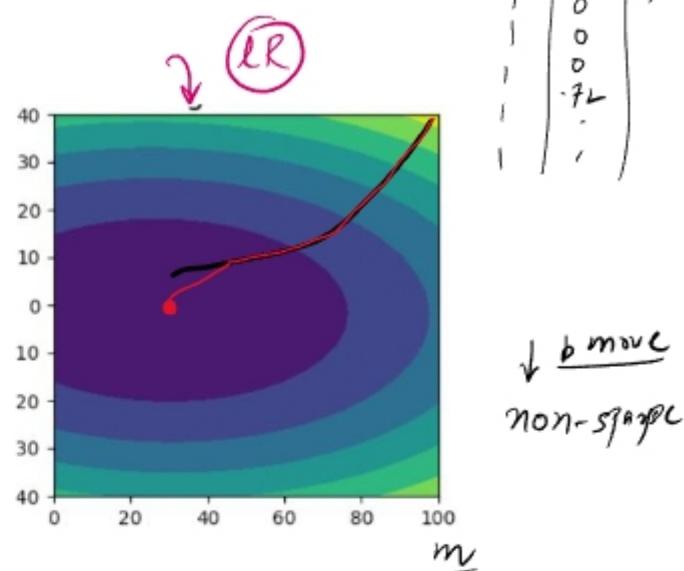
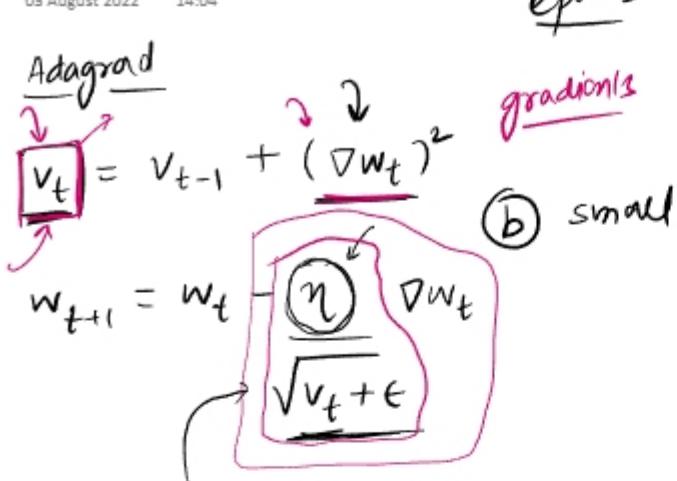


$$\begin{array}{c|c|c} x_1 & x_2 & y \\ \hline 0 & 0 & \\ 0 & 0 & \\ 1 & 0 & \\ 0 & \end{array}$$

↓
Update → small

Mathematical Formulation

09 August 2022 14:04



RMSProp

$$v_t = \beta v_{t-1} + (1-\beta)(\nabla w_t)^2$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \nabla w_t$$

exp decaying avg
old epoch

$$\beta = 0.95$$

$$\rightarrow v_0 = 0$$

$$\rightarrow v_1 = 0.95 \times 0 + 0.05 (\nabla w_1)^2$$

$$\rightarrow v_2 = \frac{0.95 \times 0.05 (\nabla w_1)^2 + 0.05 (\nabla w_2)^2}{2}$$

$$\rightarrow v_3 = \frac{0.95 \times 0.95 \times 0.05 (\nabla w_1)^2 + 0.95 \times 0.05 (\nabla w_2)^2 + 0.05 (\nabla w_3)^2}{3}$$

↑ smaller ↑ epoch 1 ↑ epoch 2 ↑ epoch 3

v_t shoot $\rightarrow \eta$ small

Adagrad \rightarrow NN \rightarrow non-convex optimization

convex optimis \rightarrow linear

$\sigma \text{ sig}$

✓ convex optimis \rightarrow linear
reg

- global minima RMSprop

Disadvantage

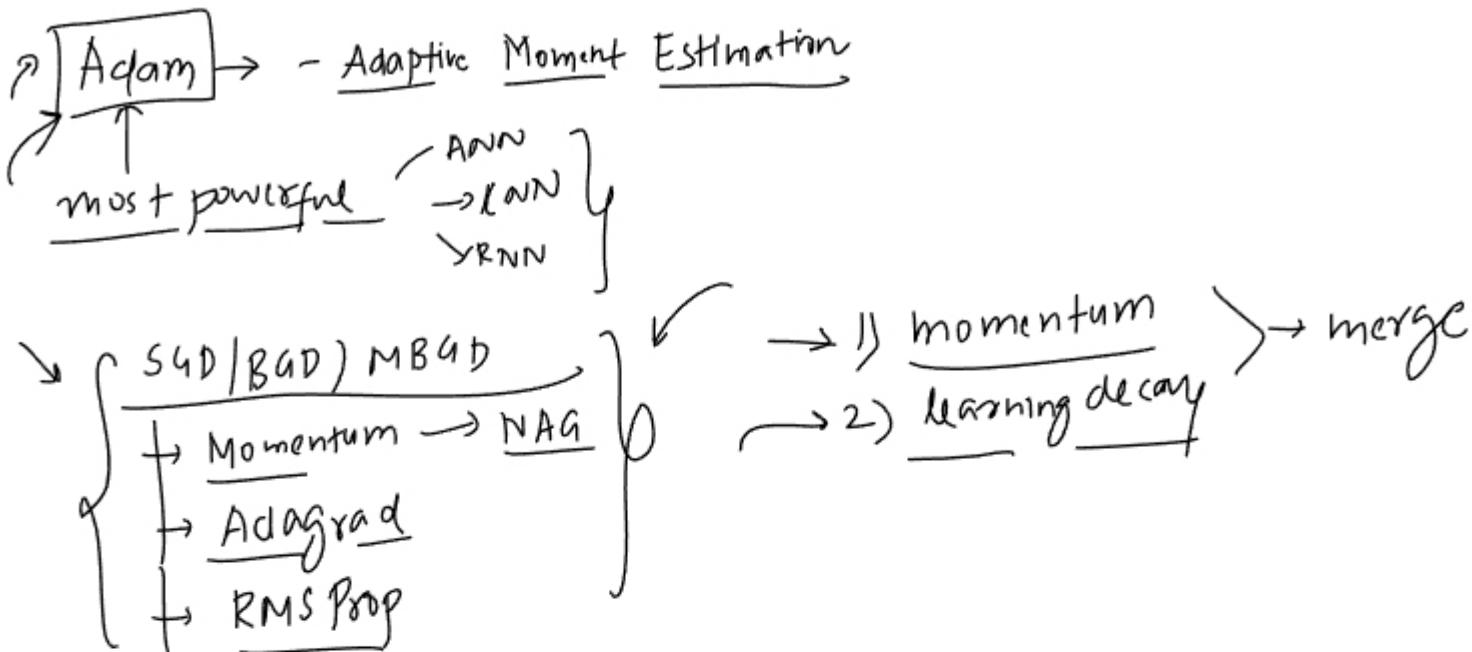
03 August 2022 14:54

No,
Optimization

↓
Adam

Introduction

04 August 2022 10:45



Mathematical Formulation

04 August 2022 10:45

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} * m_t \rightarrow \hat{m}_t = \frac{m_t}{1 - \beta_1^{[t]}} \rightarrow \hat{v}_t = \frac{v_t}{1 - \beta_2^{[t]}}$$

Keras Bias correction epoch #

where $m_t = \beta_1 m_{t-1} + (1 - \beta_1) \nabla w_t$ → momentum

$v_t = \beta_2 v_{t-1} + (1 - \beta_2) (\nabla w_t)^2$ → Adagrad

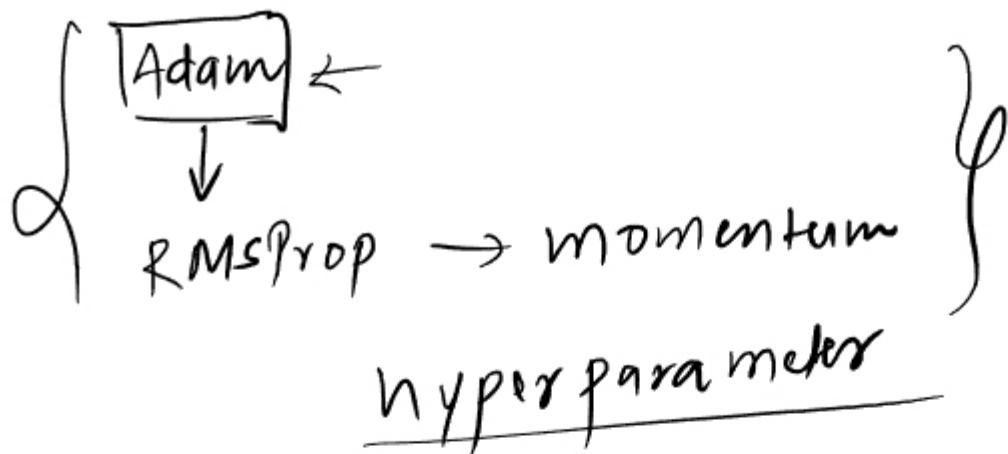
$m_0 = 0 \quad v_0 = 0 \quad \eta = 0.1 \quad \text{DO}$

Demo

04 August 2022 10:45

Verdict

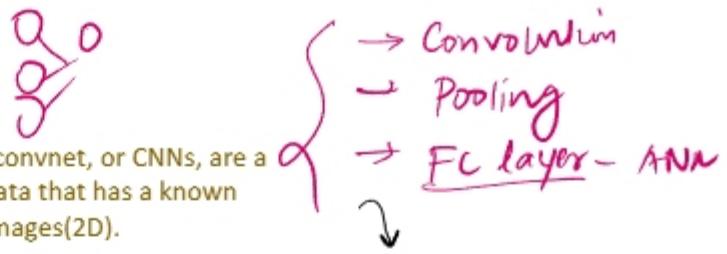
04 August 2022 10:45



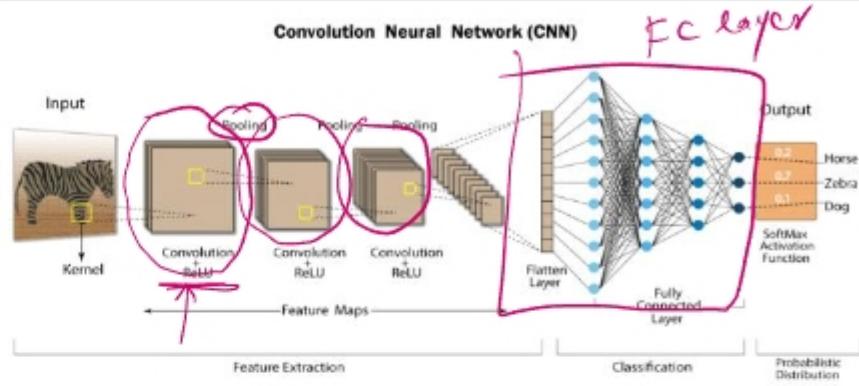
What is a CNN?

17 August 2022 06:47

Convolutional neural networks, also known as convnet, or CNNs, are a special kind of neural network for processing data that has a known grid-like topology like time series data(1D) or images(2D).



$\rightarrow \text{ANN} \rightarrow \text{Convolution}$
 $\text{CNN} \rightarrow \text{convolution} \otimes$
 $\text{ANN} \rightarrow \text{matrix mult.}$



Inspirations

\hookrightarrow visual cortex

1998 Yann LeCun \rightarrow AT&T \rightarrow
 \hookrightarrow Microsoft \rightarrow OCR hand writing $\left. \begin{array}{l} \text{recog} \\ \text{CNN} \end{array} \right\}$

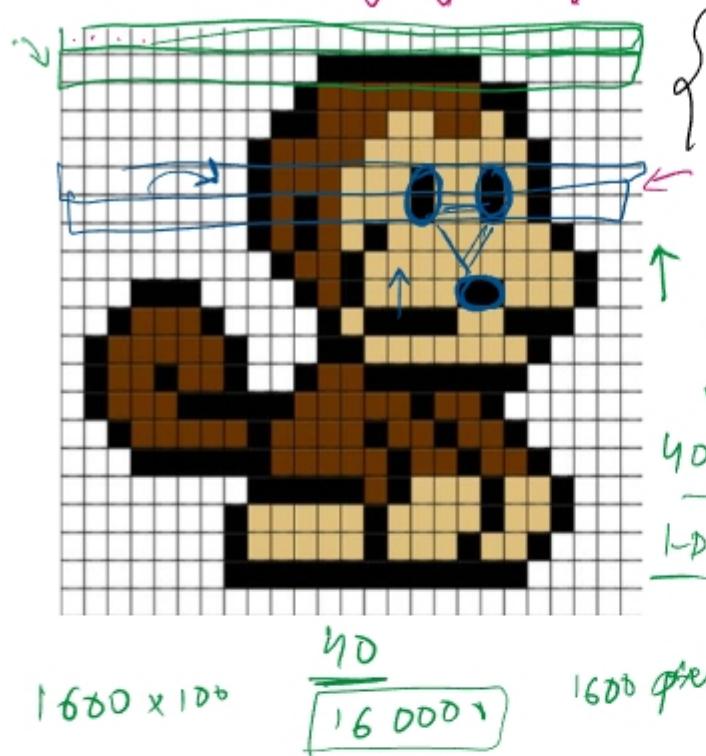
$\left. \begin{array}{l} \text{facial recos} \\ \text{self driving} \end{array} \right\} \begin{array}{c} \text{CNN} \\ \text{RNN} \end{array}$

Why not use ANN?

17 August 2022 06:47

pixels 2D grid

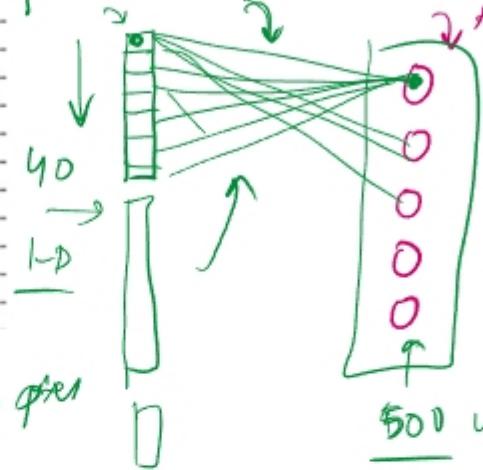
CNN vs ANN



1. High Computation Cost ✓
2. Overfitting →
3. Loss of imp info like spatial arrangement of pixels

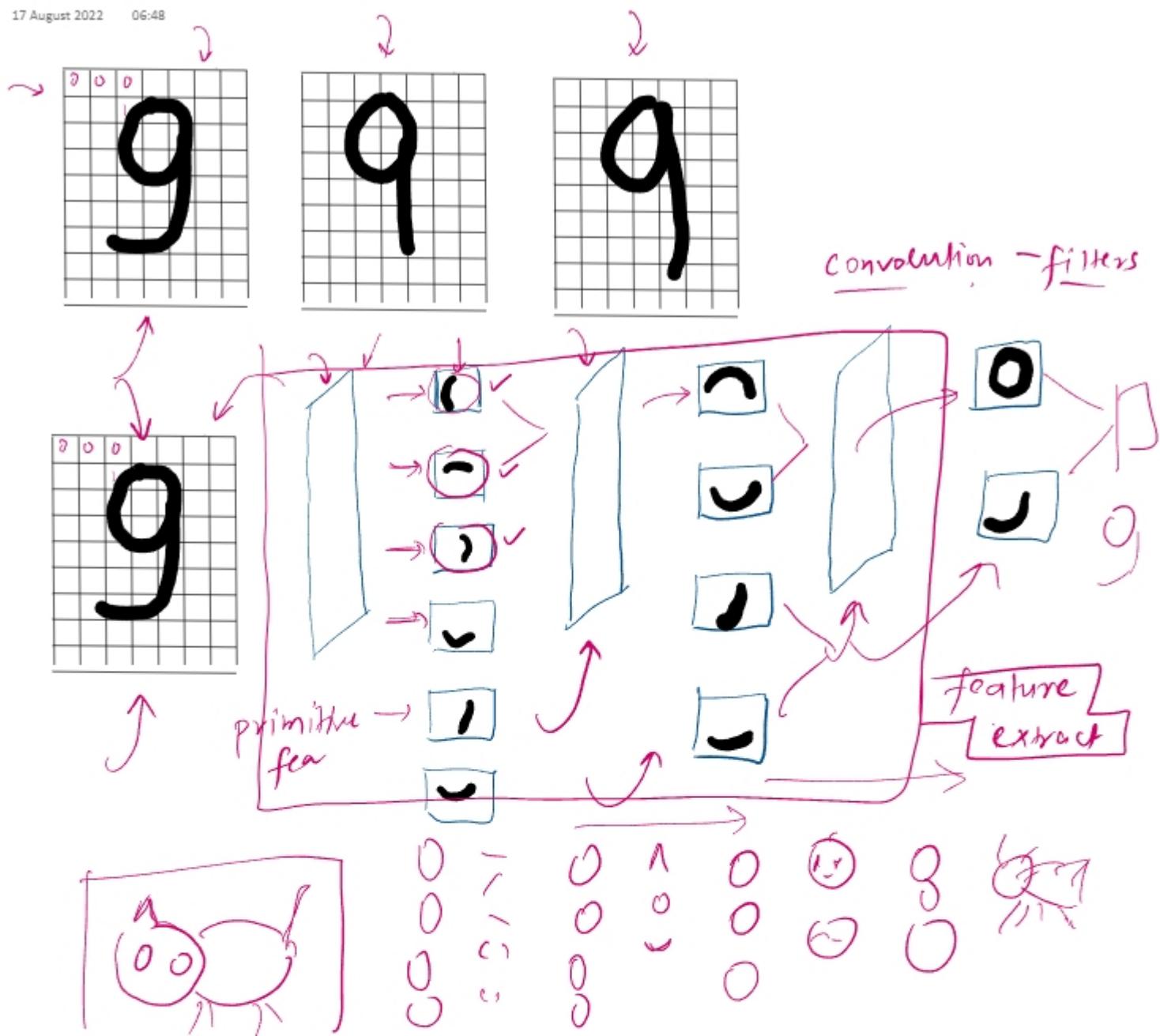
40×40 image

MNIST \rightarrow ANN \rightarrow 98%



CNN Intuition

17 August 2022 06:48

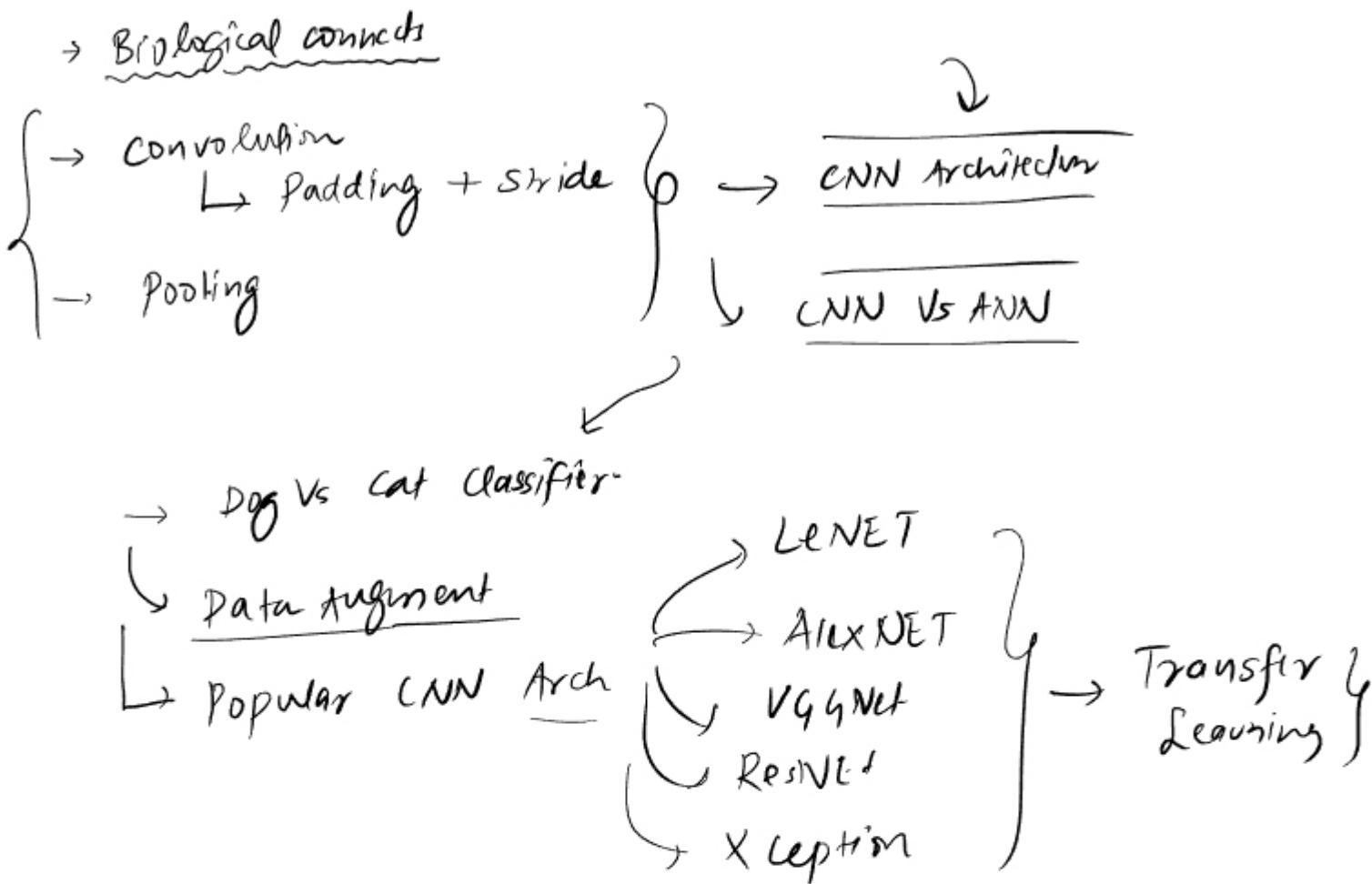


CNN Applications

17 August 2022 06:48

Roadmap

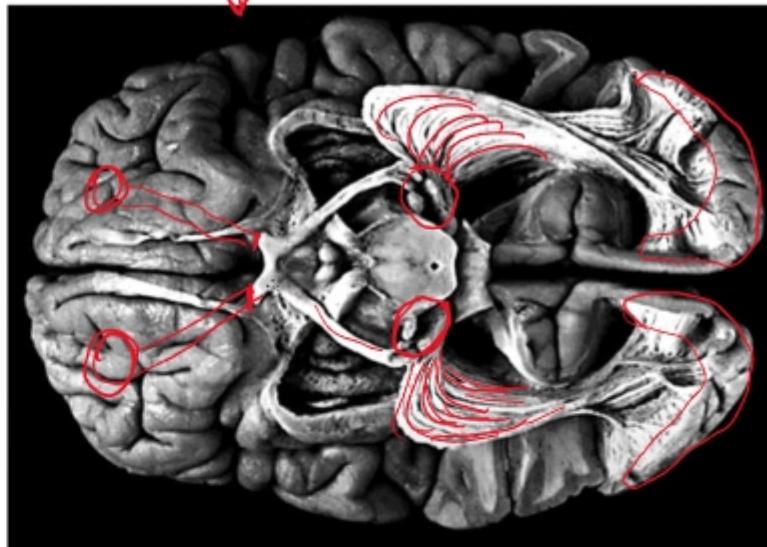
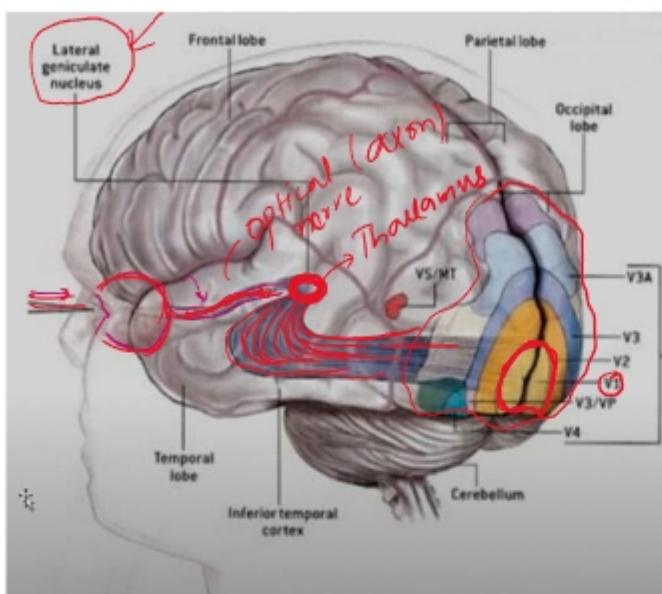
17 August 2022 06:48



Human Visual Cortex

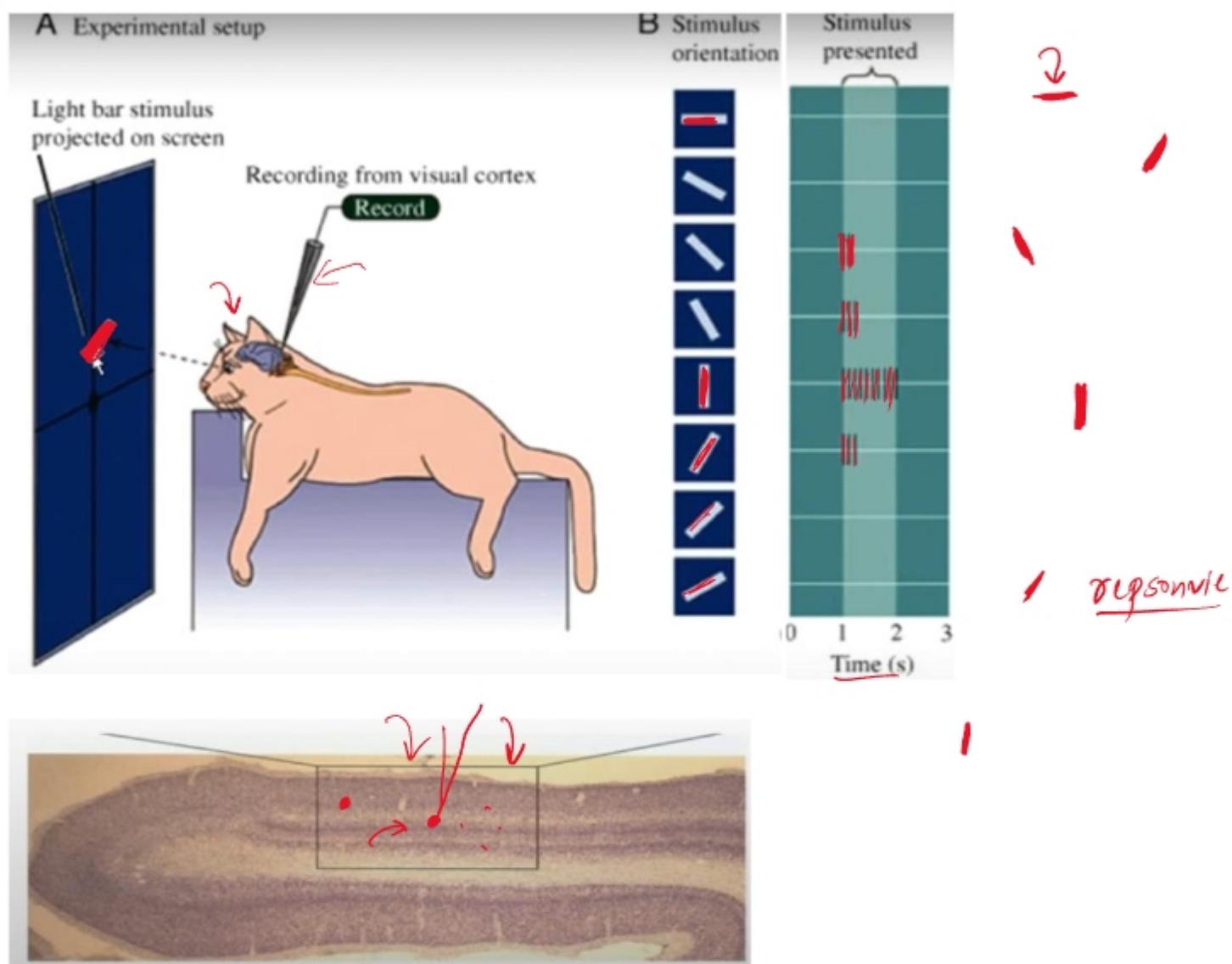
18 August 2022 16:05

28 sheet primary V cortex
actual brain



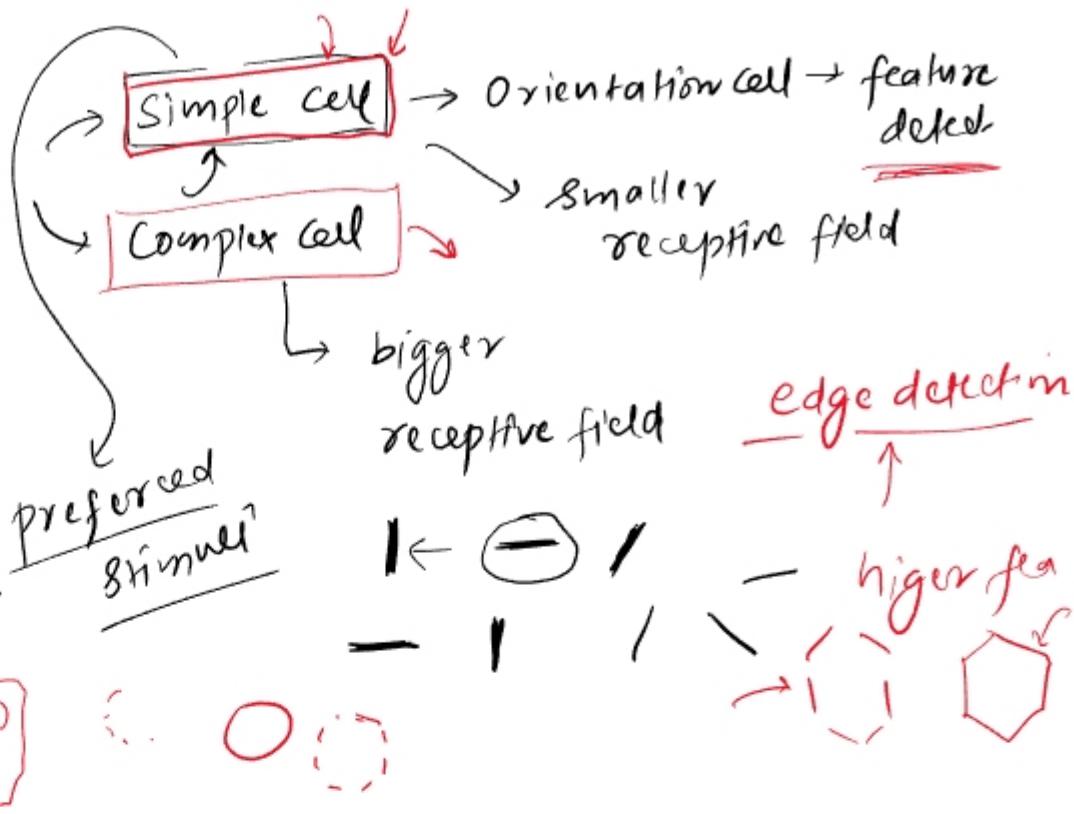
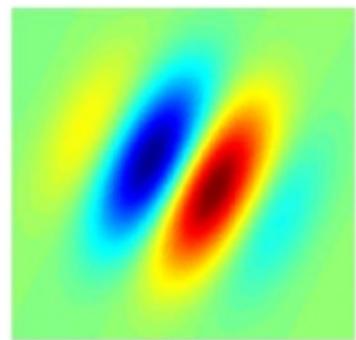
The Experiment

18 August 2022 16:09



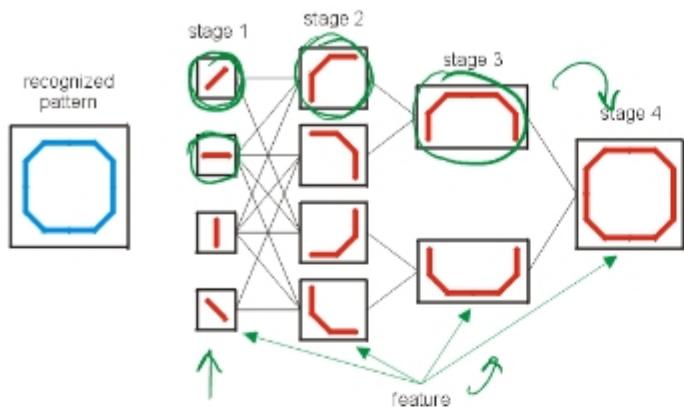
Conclusion

18 August 2022 17:57

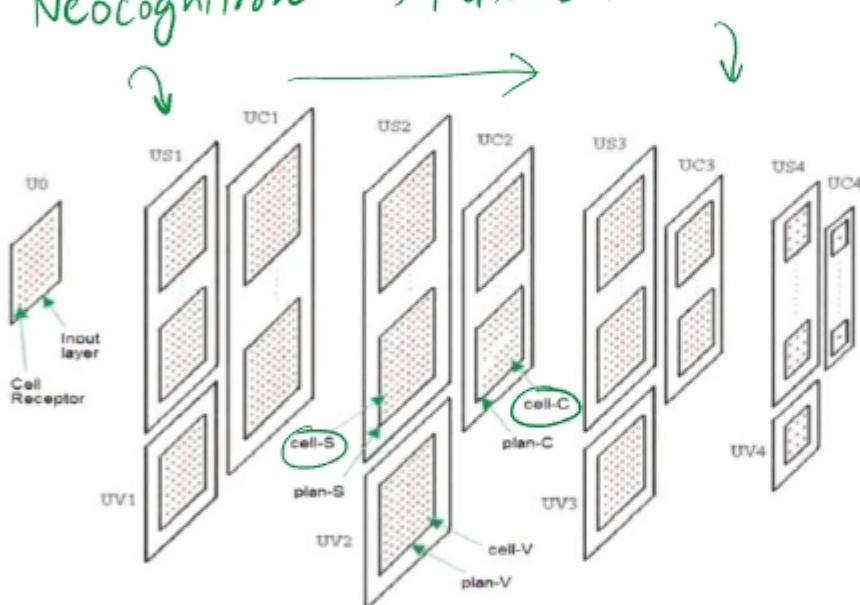


Development

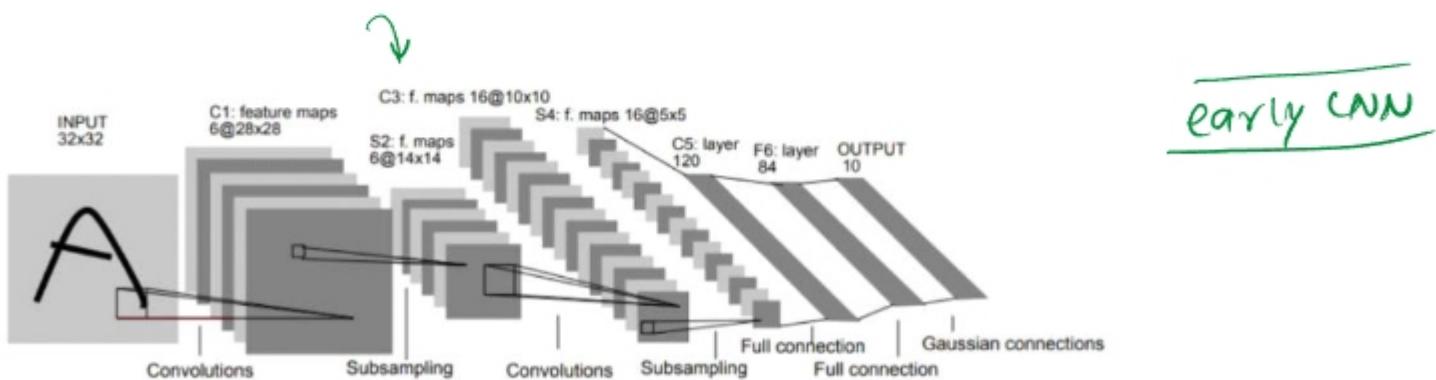
18 August 2022 16:15



Neocognitron → Fukushima



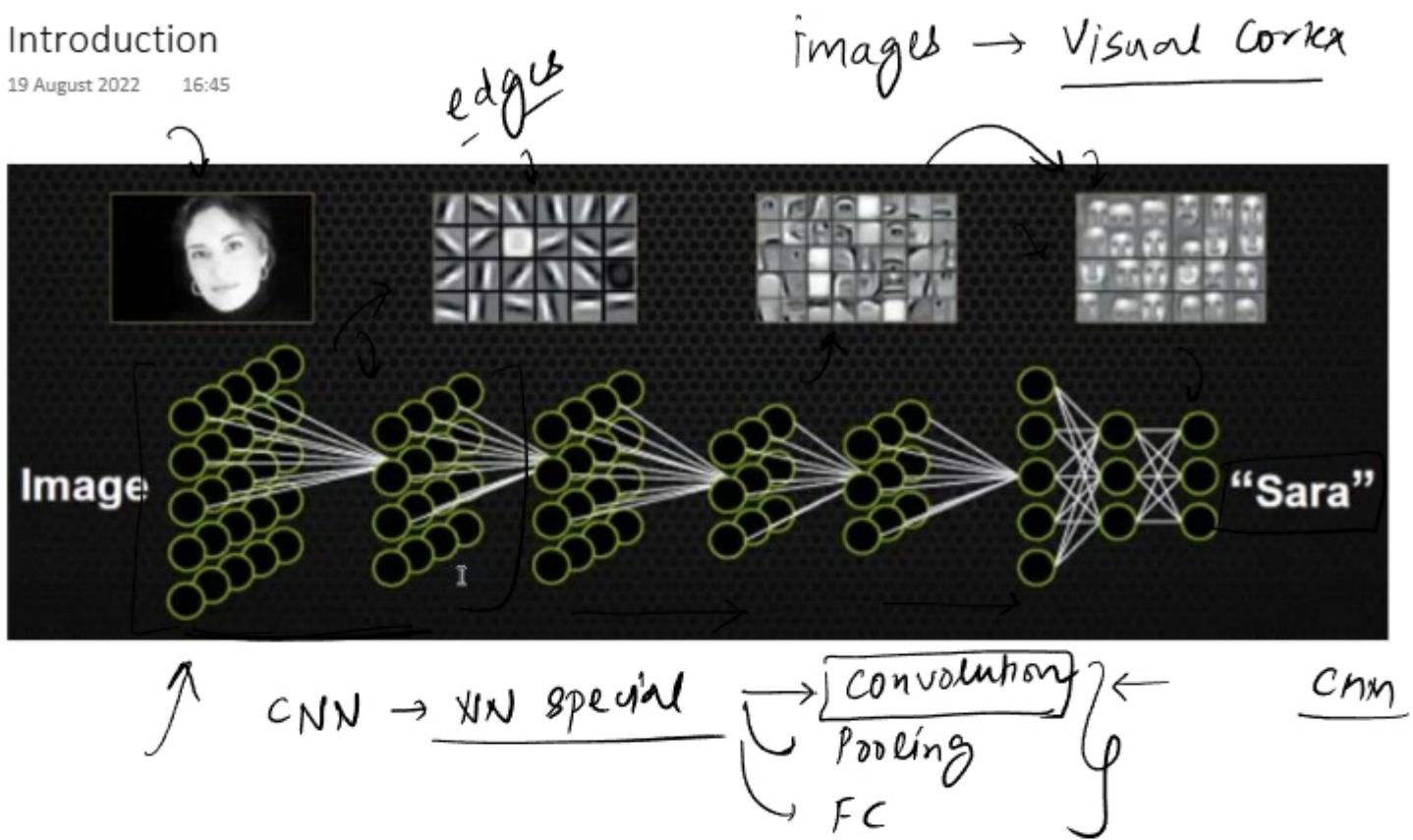
Yann LeCun → CNN → Backprop convolution



2012 → AlexNET → ImageNET → CNNs

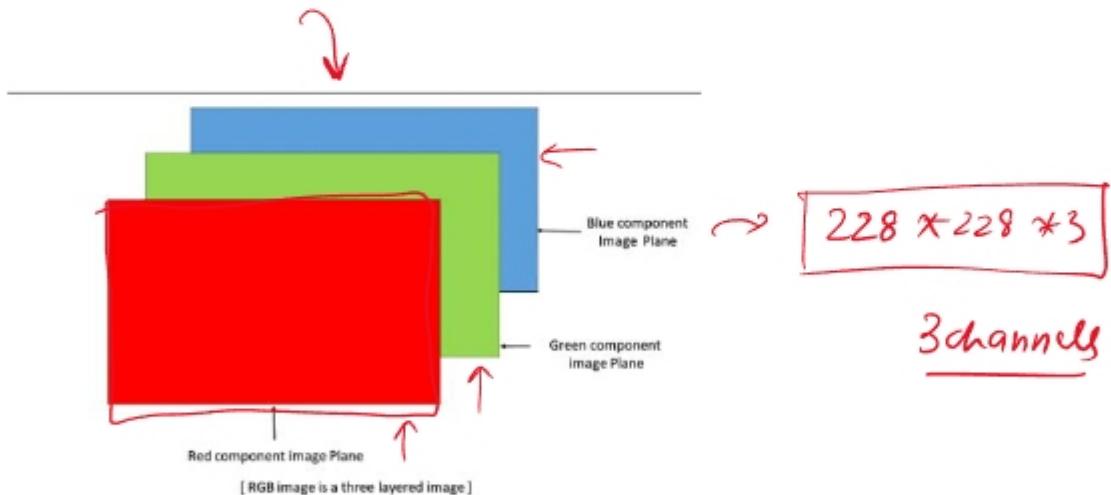
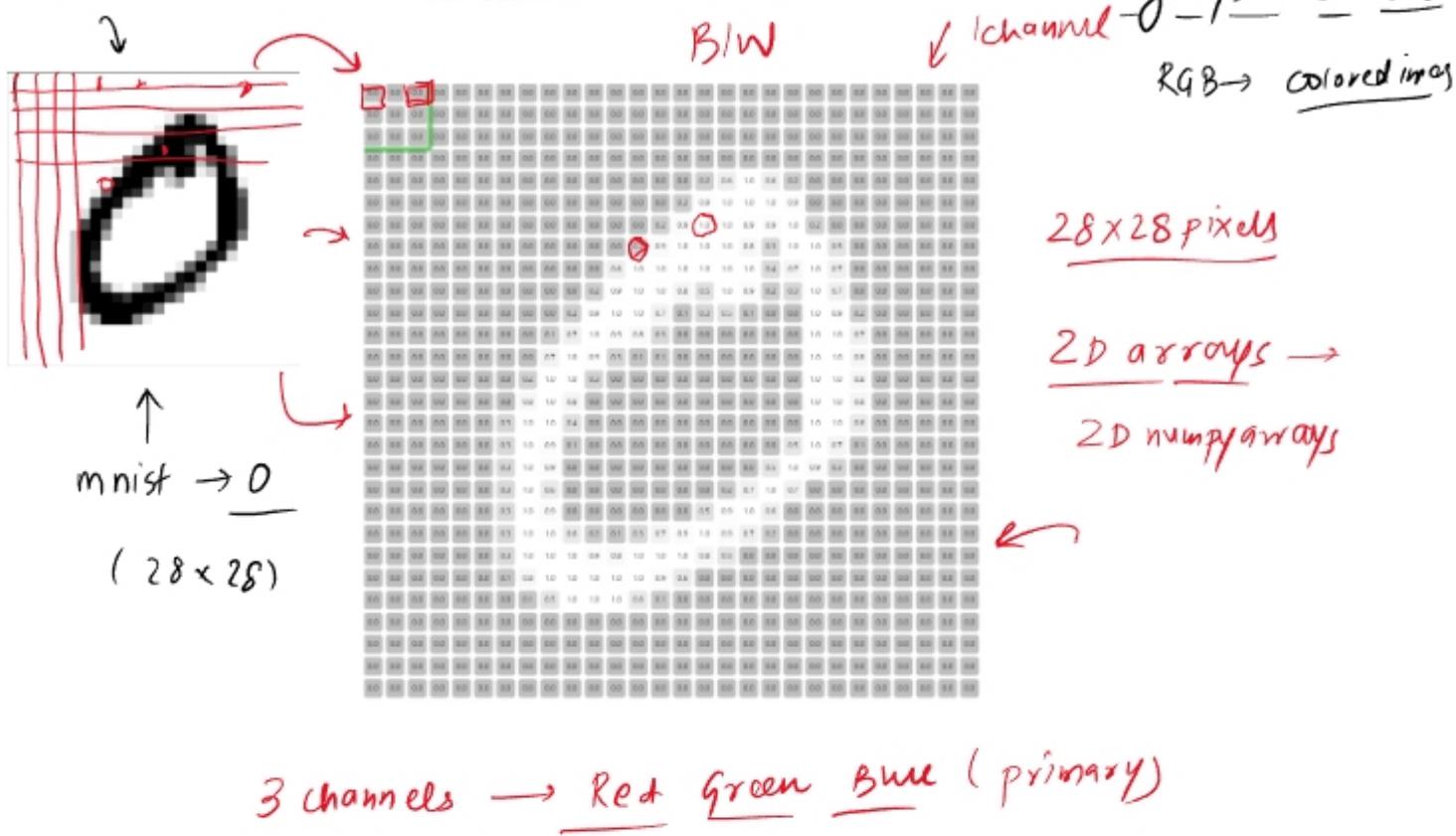
Introduction

19 August 2022 16:45



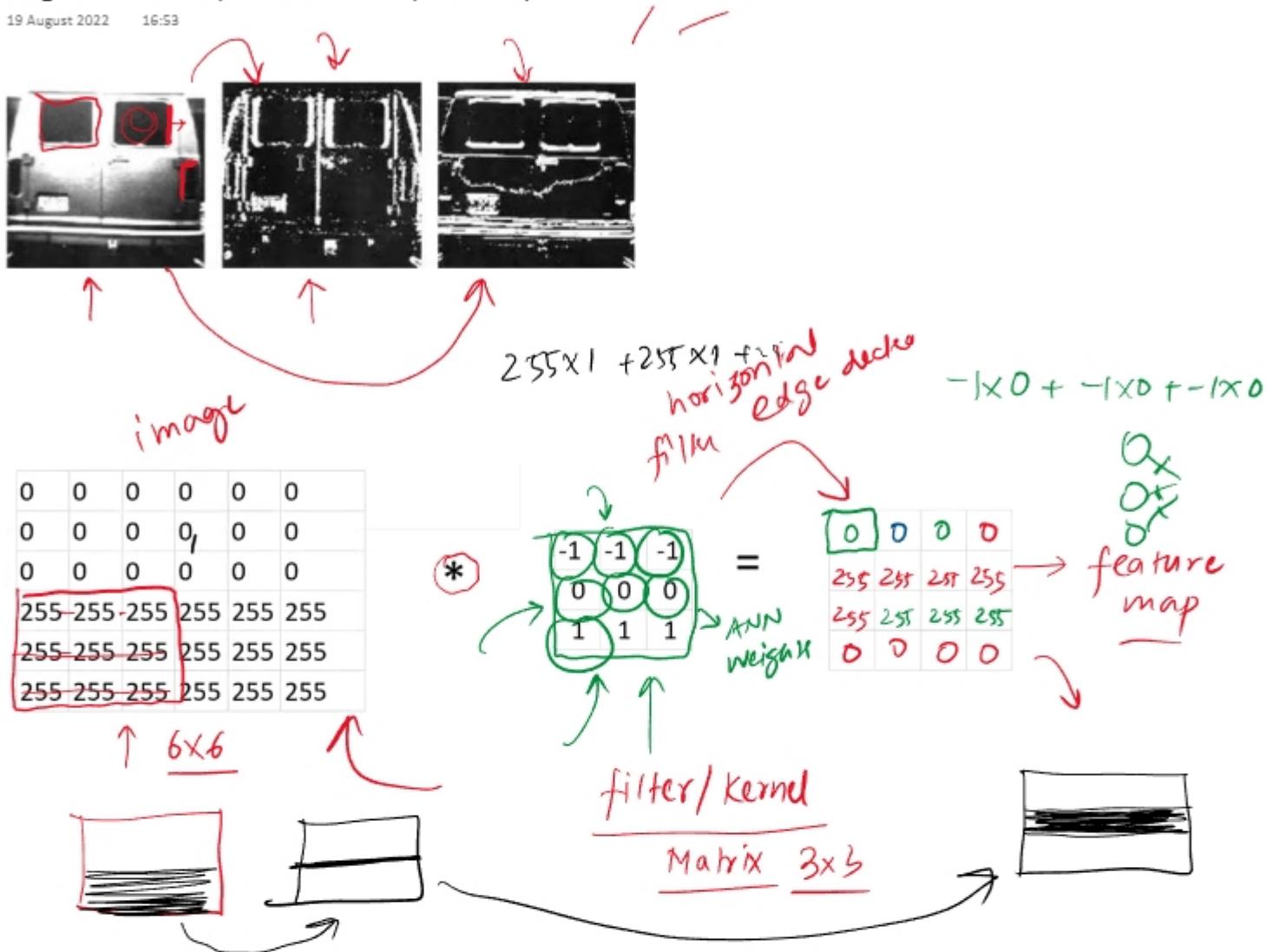
Basics of Images

19 August 2022 16:53



Edge Detection (Convolution Operation)

19 August 2022 16:53



$$\begin{matrix}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 255 & 255 & 255 & 255 & 255 & 255 \\
 255 & 255 & 255 & 255 & 255 & 255 \\
 255 & 255 & 255 & 255 & 255 & 255
 \end{matrix} * \begin{matrix}
 -1 & -1 & -1 \\
 0 & 0 & 0 \\
 1 & 1 & 1
 \end{matrix} = \begin{matrix}
 & & & & & \\
 & & & & & \\
 & & & & & \\
 & & & & & \\
 & & & & & \\
 & & & & &
 \end{matrix}$$

(6×6) (3×3) (4×4)

$$\begin{matrix}
 (28 \times 28) \rightarrow (26 \times 26) & (3 \times 3) \rightarrow ? \underline{(26 \times 26)} \\
 n \times n & m \times m \rightarrow (n-m+1) \times (n-m+1)
 \end{matrix}$$

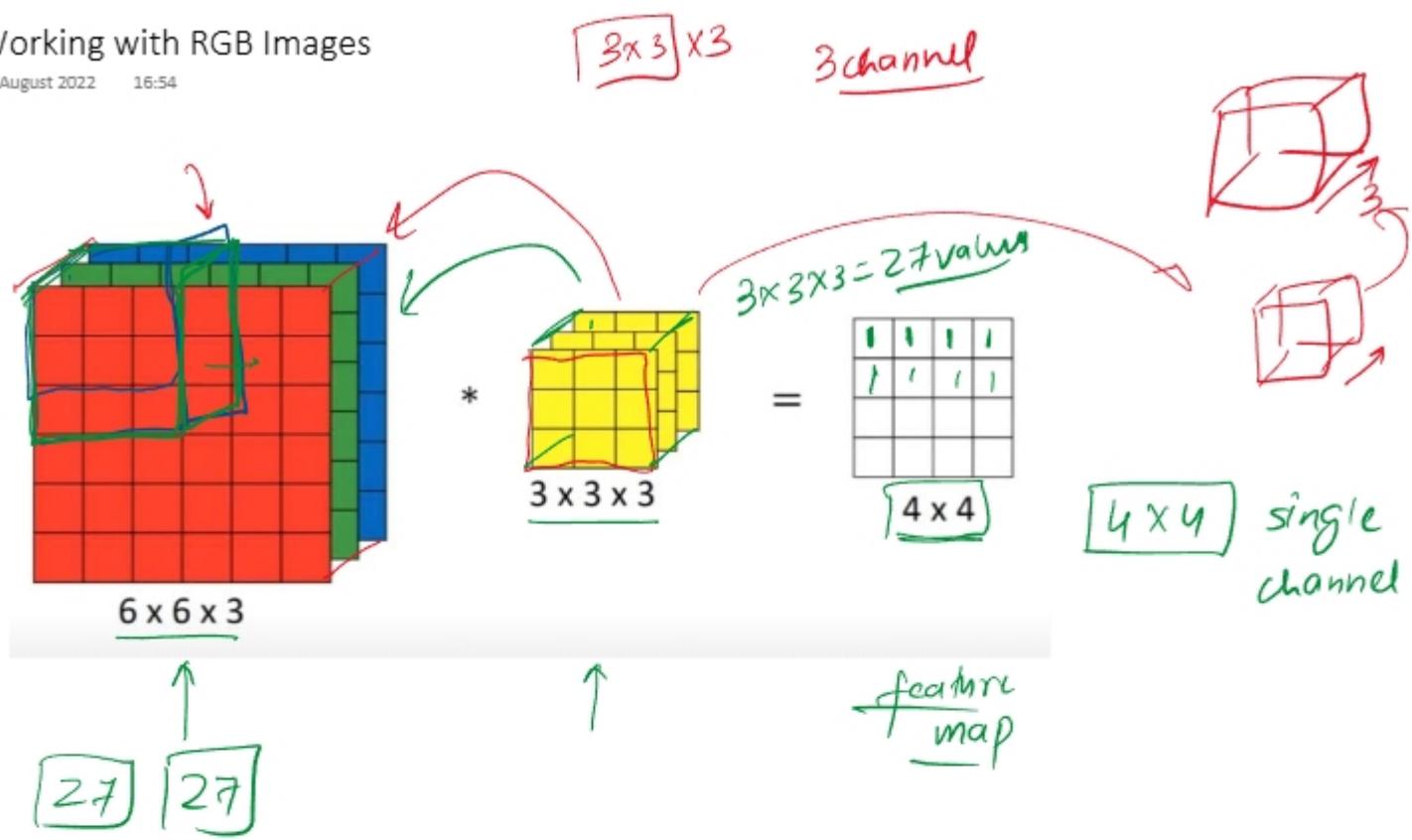
$$(64 \times 64) \quad (3 \times 3) \rightarrow \underline{(62 \times 62)}$$

Demo

19 August 2022 16:54

Working with RGB Images

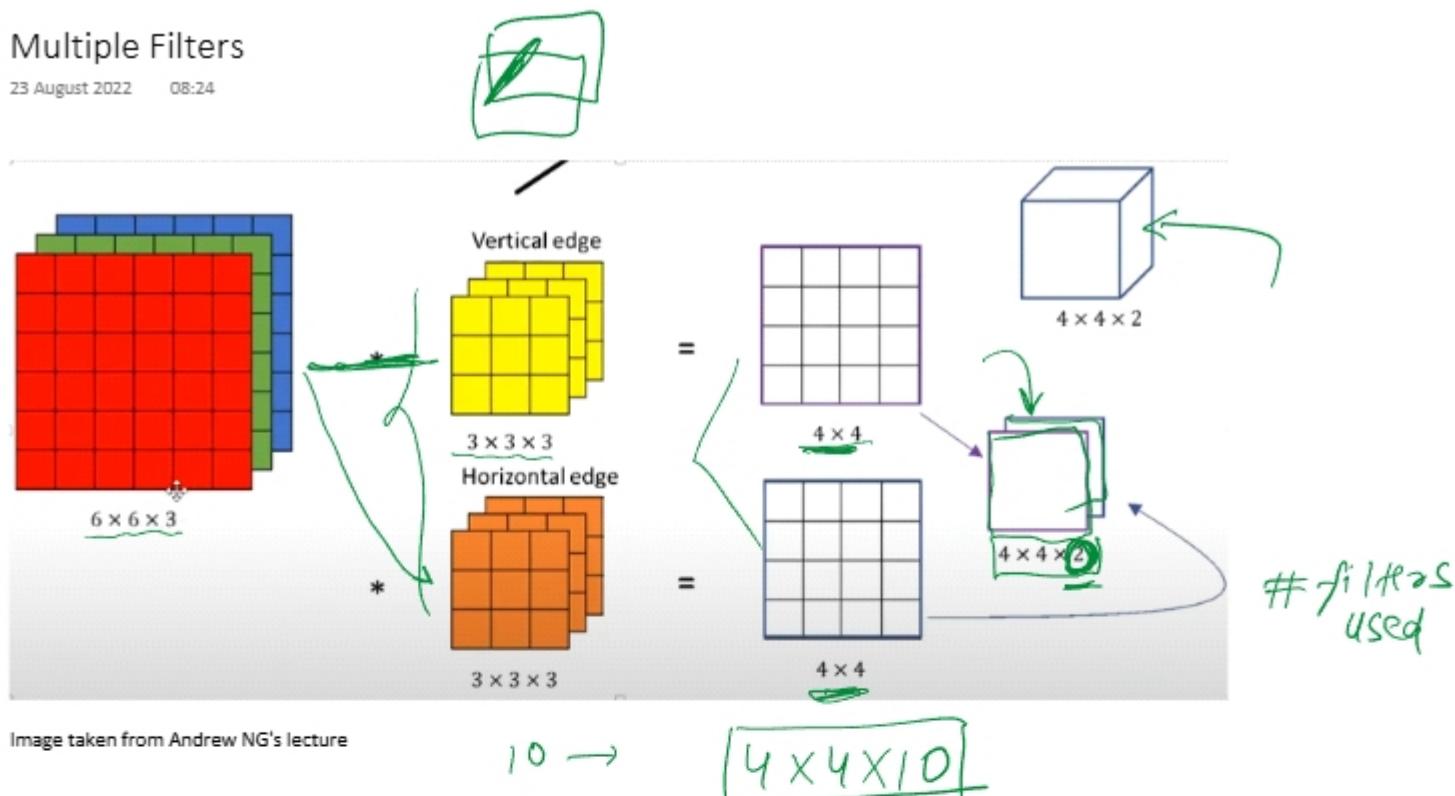
19 August 2022 16:54



$$[m \times m \times C] \quad [n \times n \times C] \rightarrow \frac{(m-n+1) \times (m-n+1)}{\text{single channel}}$$

Multiple Filters

23 August 2022 08:24



Problem with Convolution

26 August 2022 14:25

What is Padding?

26 August 2022 14:26

Handwritten annotations:

- $\eta \times n$
- 5×5
- 7×7
- $f \times f$
- 3×3
- $7 \times 1 + 4 \times 1 + 3 \times 1 + 2 \times 0 + 5 \times 0 + 3 \times 0 + 3 \times -1 + 3 \times -1 + 2 \times -1 = 6$
- $n - f + 1 = n$
- $(n - f + 1) (n - f + 1)$
- $(5 - 3 + 1) = 3 \times 3$
- $n - f + 1 = 5$
- $n - 3 + 1 = 5 \Rightarrow n = 8 - 1$
- $n = 7$

0	0	0	0	0	0	0
0	7	2	3	3	8	6
0	4	5	3	8	4	0
0	3	3	2	8	4	0
0	2	8	7	2	7	0
0	5	4	4	5	4	0
0	0	0	0	0	0	0

7×7
Zero
padding

Convolution

$5 \times 5 \rightarrow 3 \times 3$

$$\underline{5 \times 5} \rightarrow 3 \times 3$$

$$(n - f + 1)$$



$$(n + 2p - f + 1)$$

$$= 5 + 2(1) - 3 + 1 = 5$$

Keras

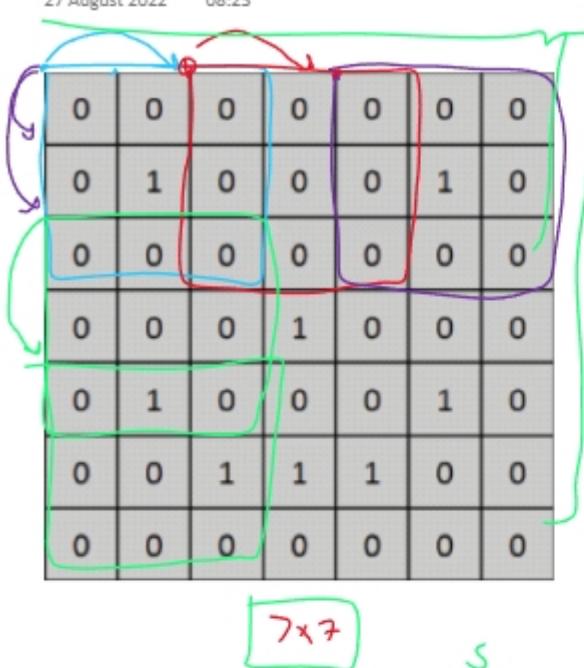
Valid

Same



Strides

27 August 2022 08:23



Stride = 1

0	0	1
1	0	0
0	1	1

3x3

(1,1)

3x3

→ right
bottom

3x3

result
feature

Stride = (2,2)

Stride = 2 →

$$\frac{7-3}{2} + 1$$

$$2+1=3$$

$$(n-f+1) \rightarrow \left[\frac{n-f}{s} + 1 \right] \rightarrow P=p$$

$$\rightarrow \left[\frac{n+2p-f}{2} + 1 \right] \rightarrow \text{stroked convolution}$$

$$\frac{7+2-3}{2} + 1 = [4 \times 4]$$

$$\frac{n-f}{2} \quad \frac{6-3}{2} \quad 1.5 = 1 + 1 = 2$$

Special Case

Stride = 2



2	3	4
---	---	---

$$\frac{6 \times 7}{3 \times 3}$$

--	--	--

$$\left\lceil \frac{n-f}{s} + 1 \right\rceil$$

$$19 \rightarrow 1 \\ 1 \cdot 1 \rightarrow 1 \\ f \text{ for } r$$

9	8	3	6	7	9	3
8	0	9	4	7	2	1
9	10	12	6	9	8	0

7x6

$$\begin{aligned}
 & \left[\frac{s}{s} + 1 \right] \rightarrow \text{floor} \\
 & \left[2 + \frac{6-3}{2} + 1 \right] = \boxed{1.5+1} \\
 & \frac{7-3}{2} + 1 = 2 + 1 = 3
 \end{aligned}$$

Why Strides are required?

27 August 2022 08:24

1) High level features

2) Computing →

Keras → stride

$$\left[\frac{n + 2P - f}{s} + 1 \right]$$

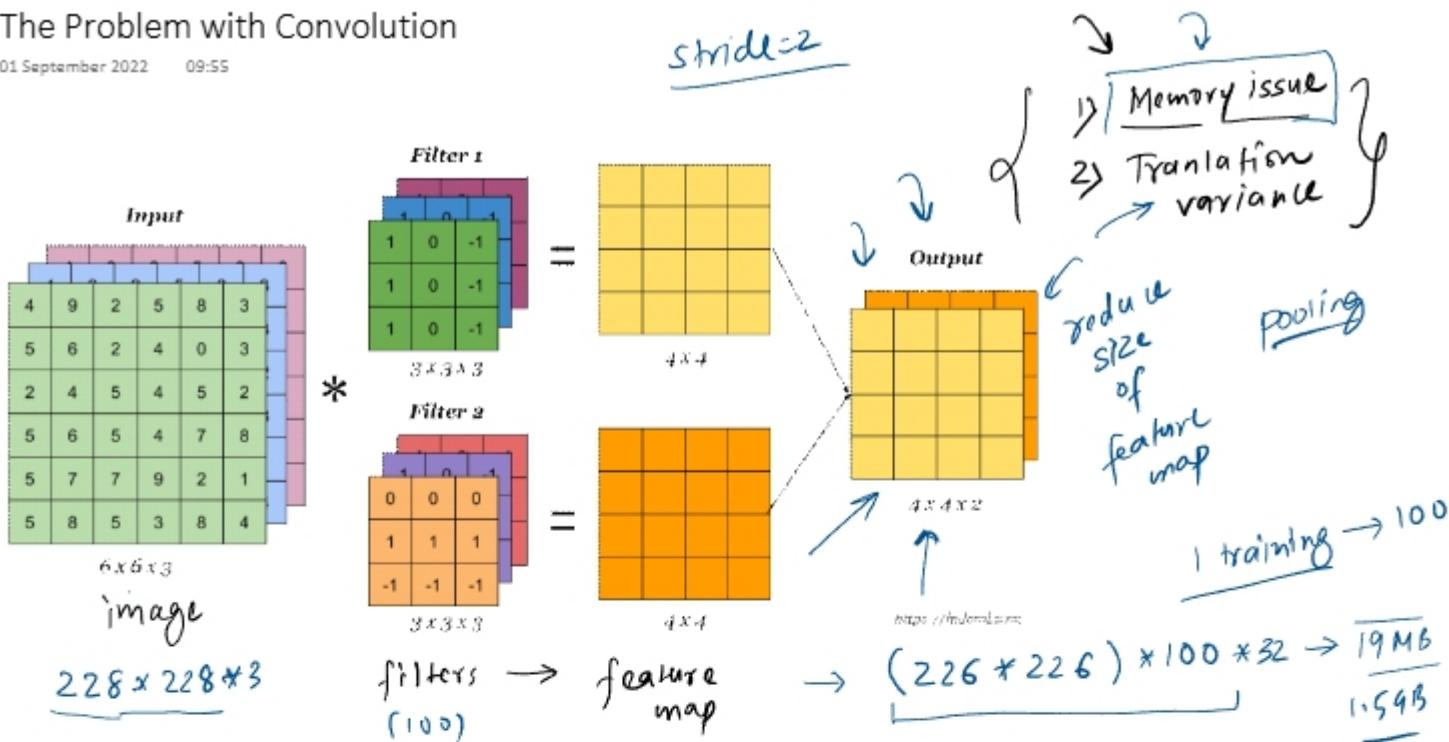
$$\frac{28 + 2 - 3}{2} + 1$$

$$\underline{13.5 + 1}$$

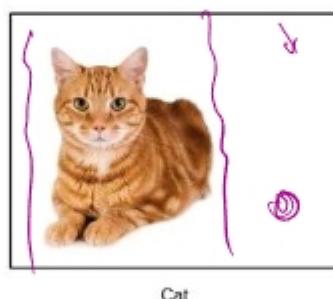
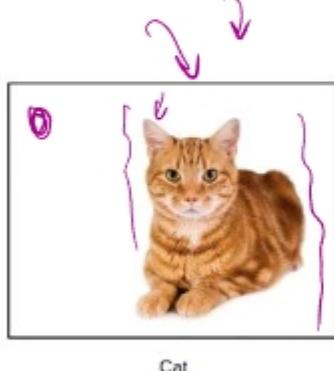
$$13 + 1 = 14$$

The Problem with Convolution

01 September 2022 09:55



Translation Variance



features
{ location dependent }

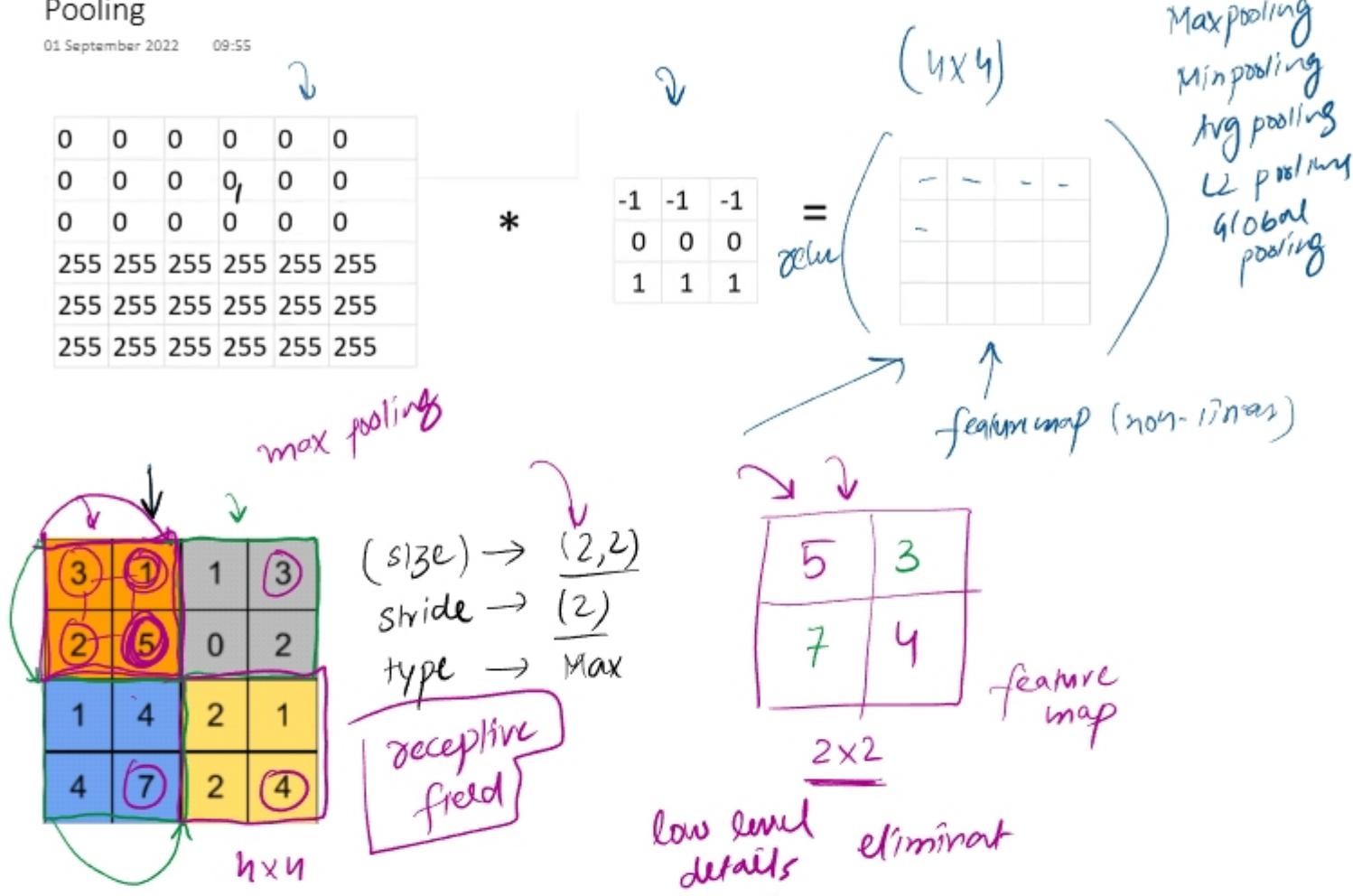
location

down sample
your
feature
map

pooling

Pooling

01 September 2022 09:55

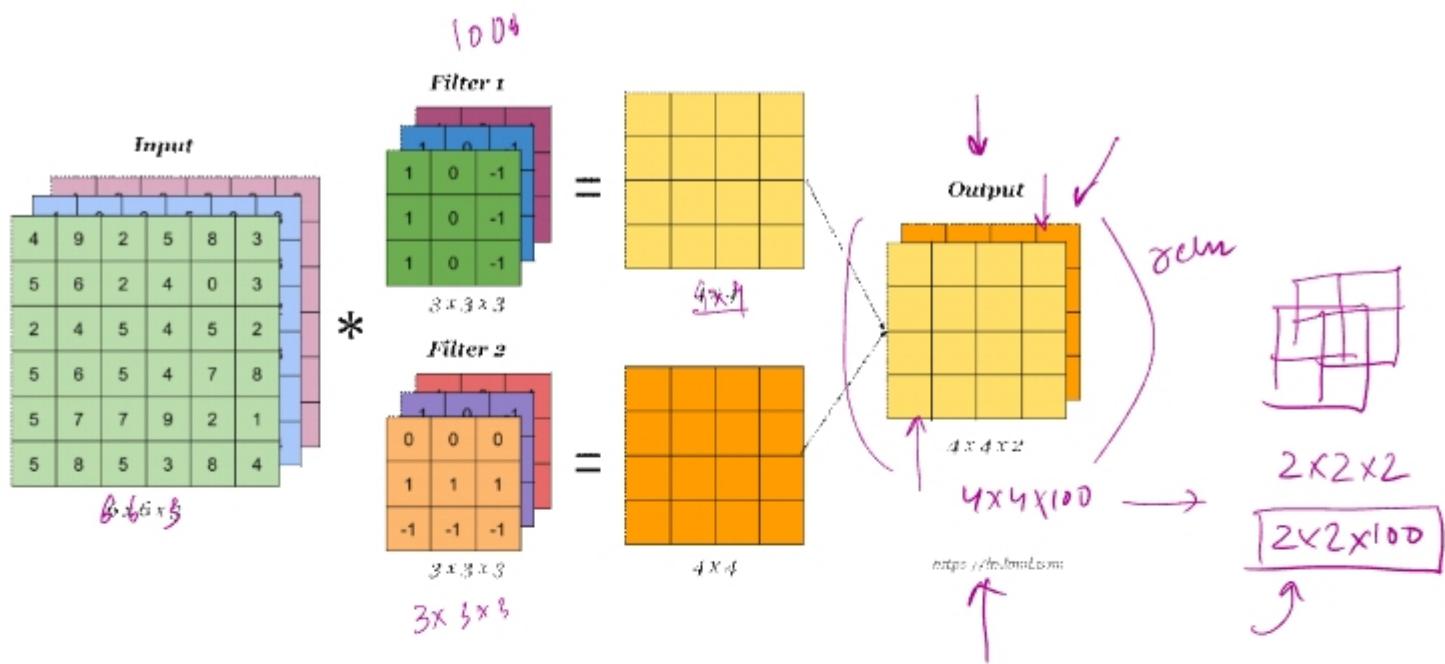


Demo

01 September 2022 09:57

Pooling on Volumes

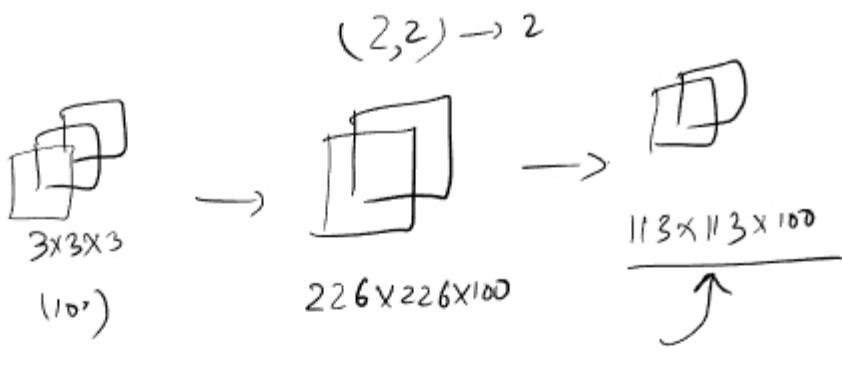
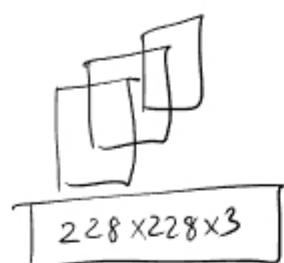
01 September 2022 09:56



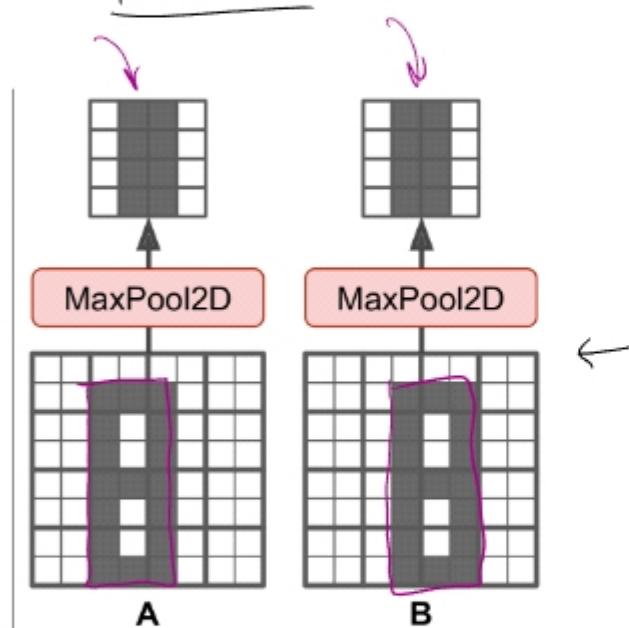
Advantages of Pooling

01 September 2022 09:56

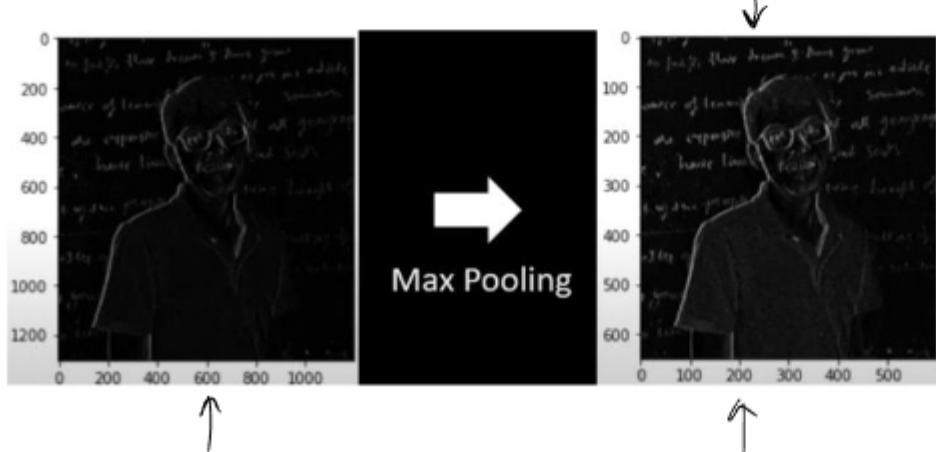
- 1) reduced size



- 2) Translation invariance



- 3) Enhanced features
(only in case of Maxpooling)



4) No need of training



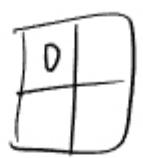
min

max pooling

avg pooling

faster

aggregate



(2×2)

2

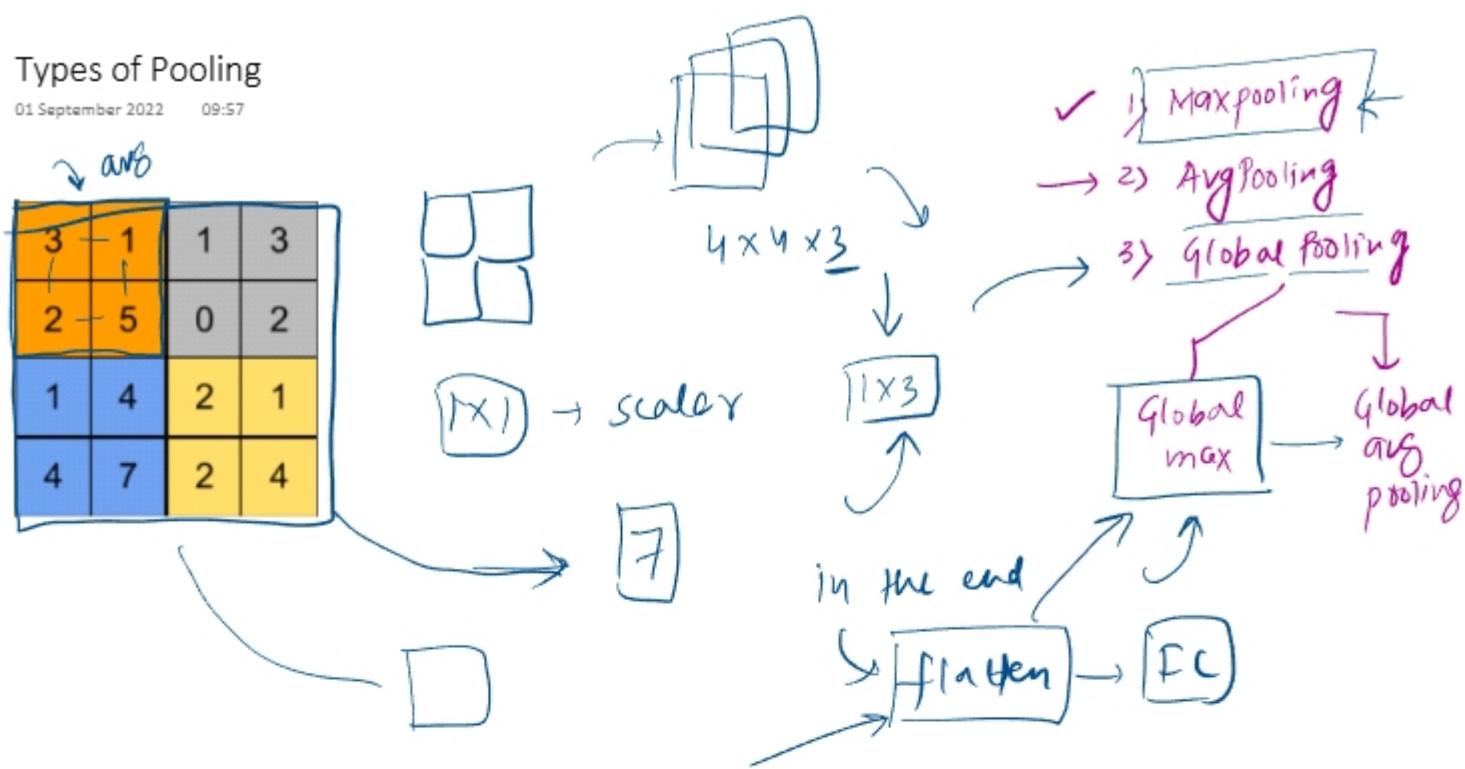
type

Keras Code

01 September 2022 09:56

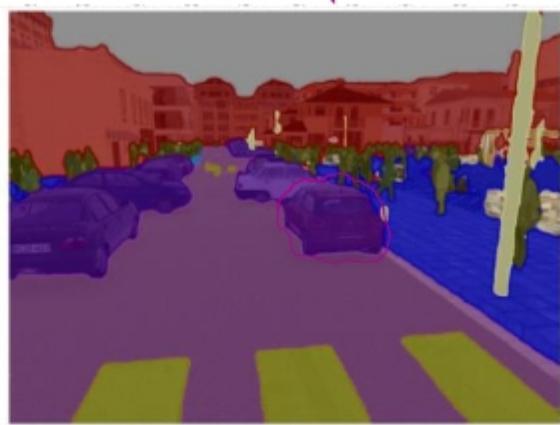
Types of Pooling

01 September 2022 09:57



Disadvantages of Pooling

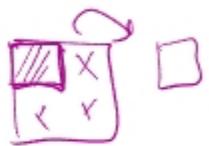
01 September 2022 09:57



lose a lot of info

$4 \times 4 \rightarrow 16$

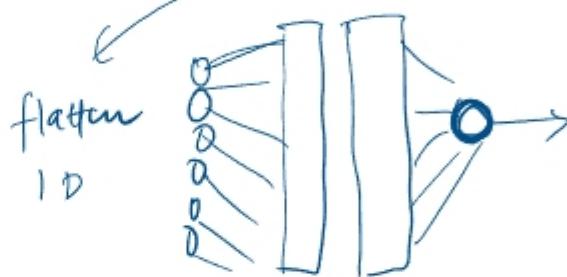
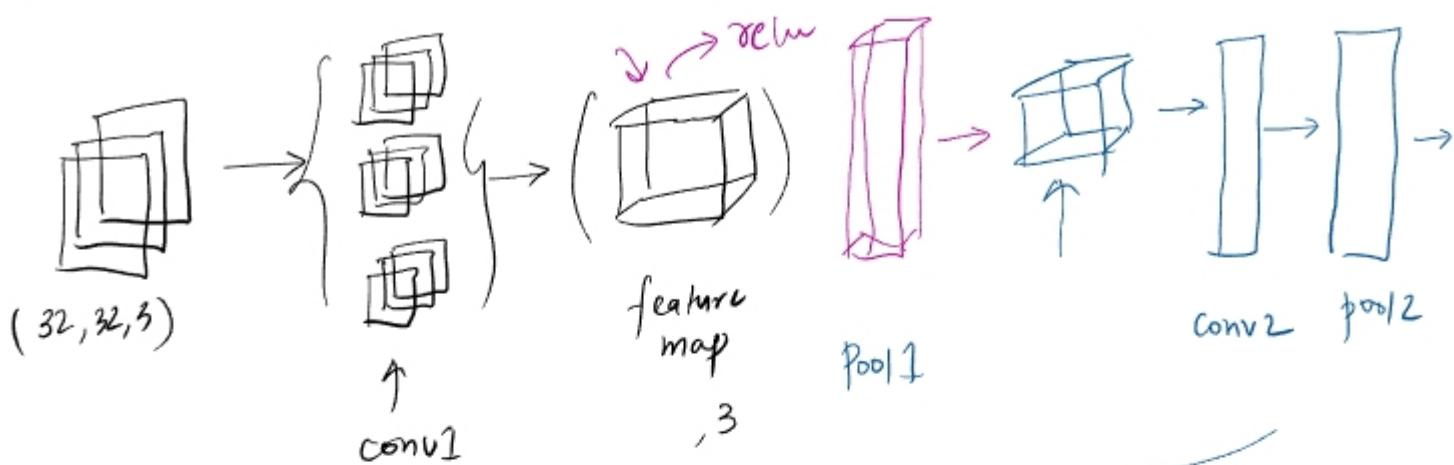
$2 \times 2 \rightarrow 4$



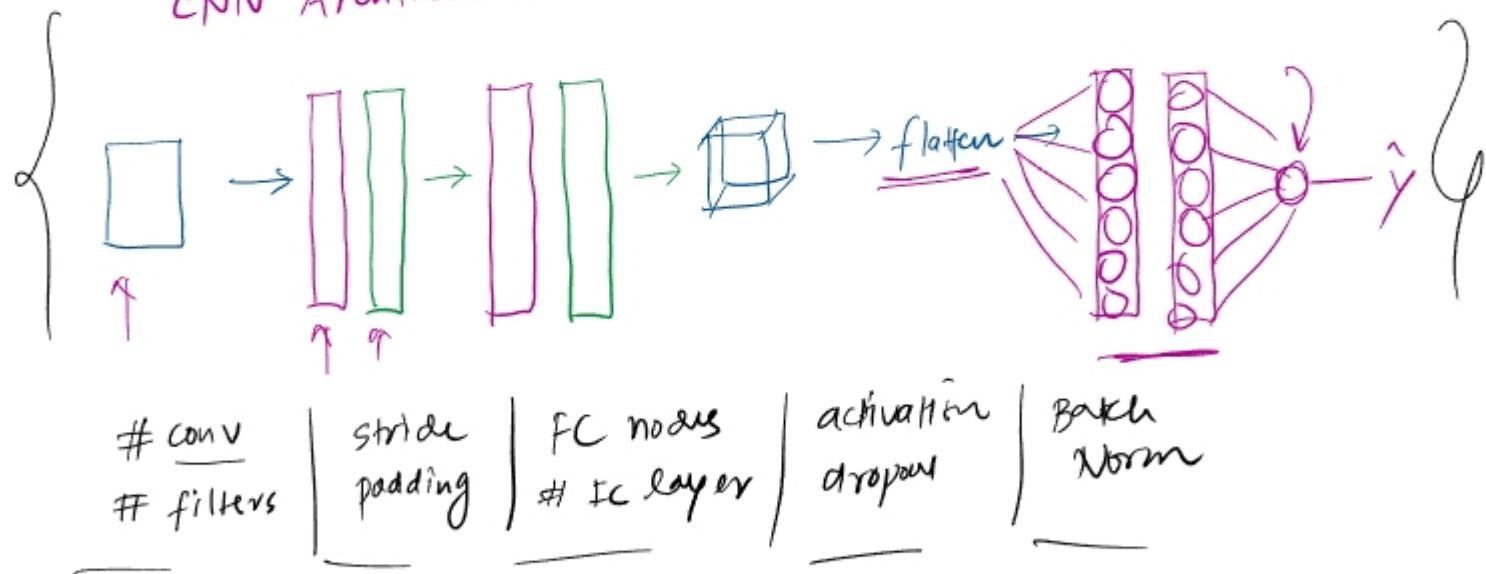
CNN Architecture

02 September 2022 10:55

1) convolution \Rightarrow Padding / stride \Rightarrow Pooling



\searrow CNN Architecture



ImageNET

1) [LeNET] \rightarrow Yann LeCUNN

2) AlexNET

3) GoogLeNET

4) VggNET

5) ResNET

6) Inception

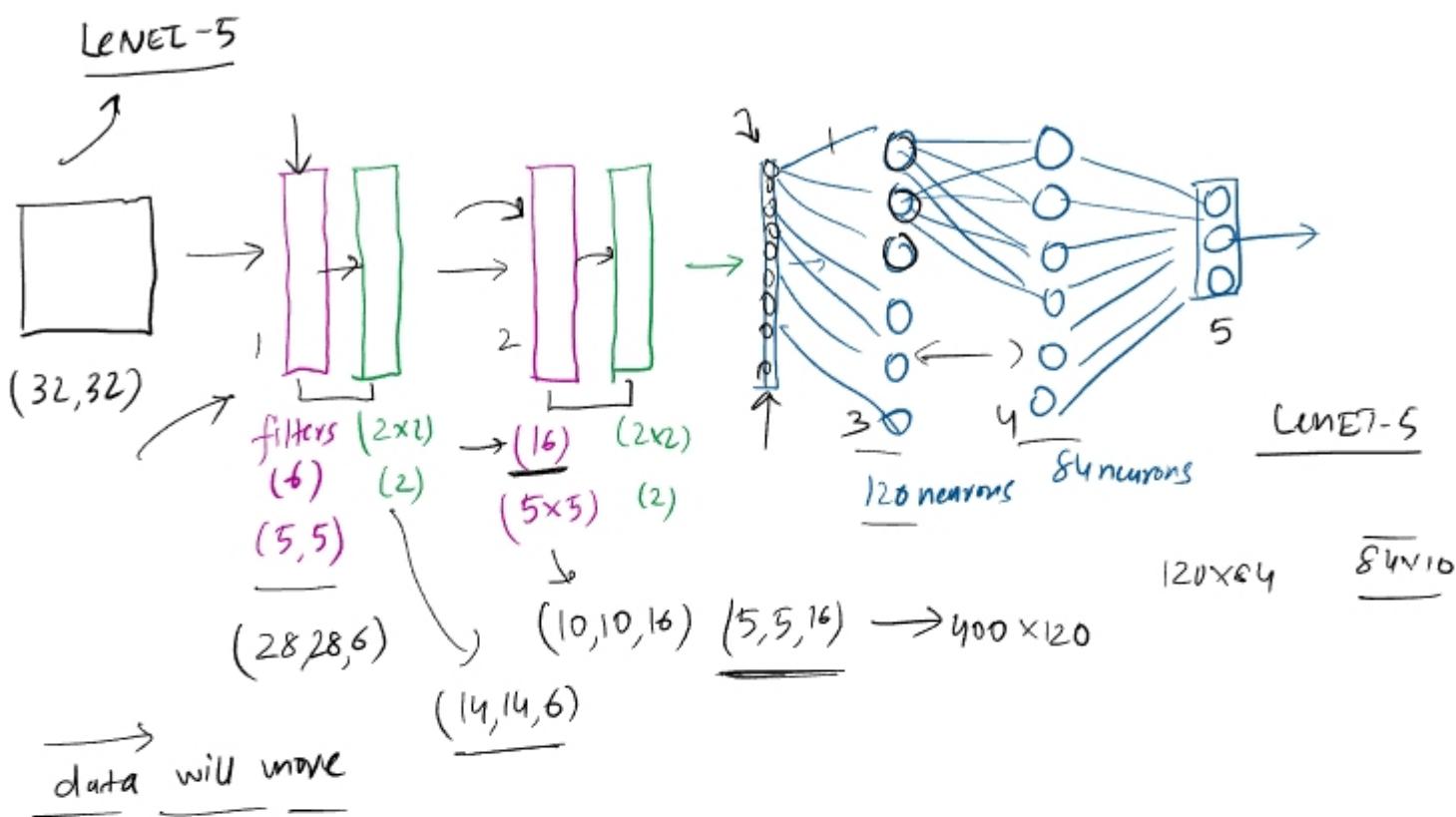
3) GoogLeNET

卷积神经网络

LeNet

02 September 2022 10:55

yann lecun \rightarrow 1989 \rightarrow 1998 \rightarrow LeNET \rightarrow CNN \rightarrow US Navy Postal Service



AlexNET
VggNET
ResNET
Inception

Guidelines

02 September 2022 10:58

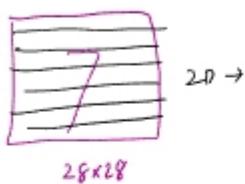
Keras Code

02 September 2022 14:58

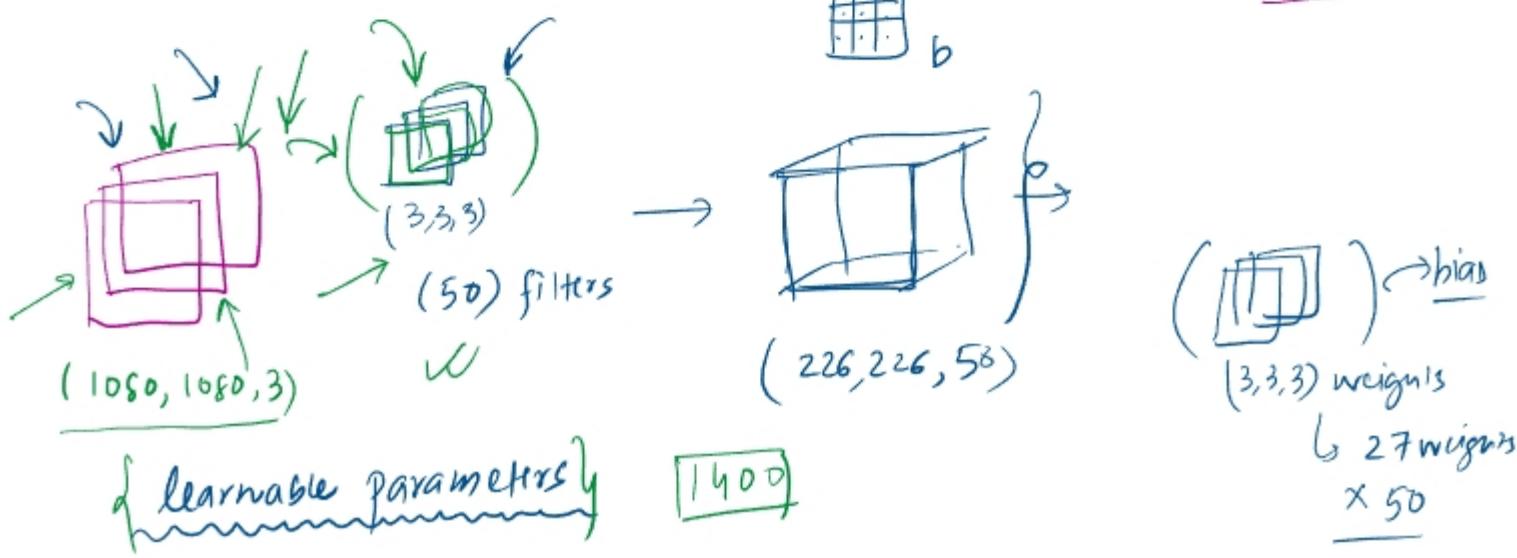
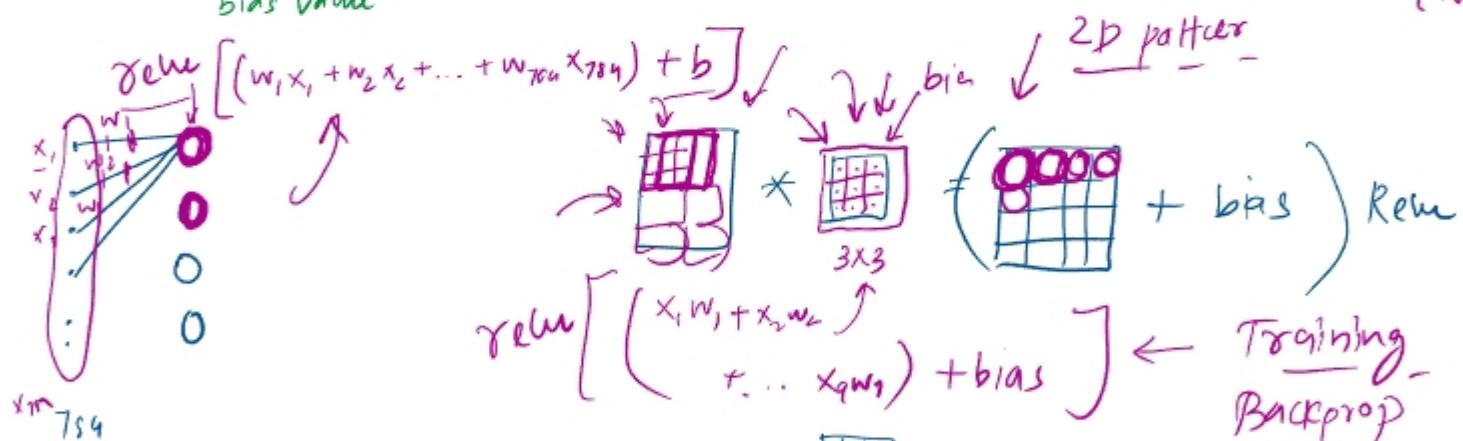
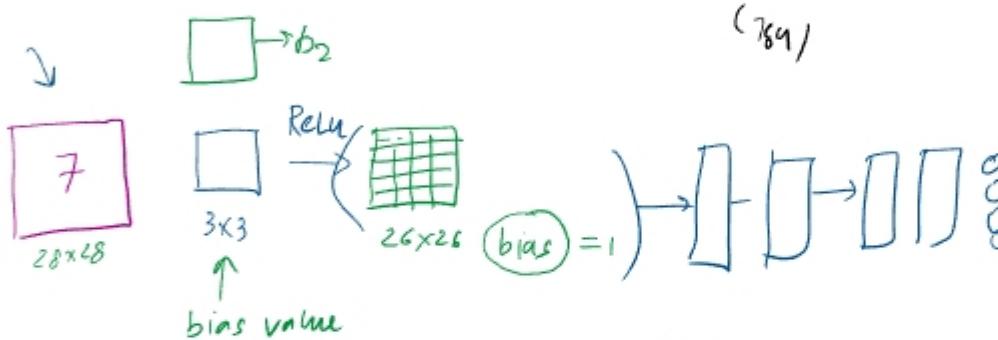
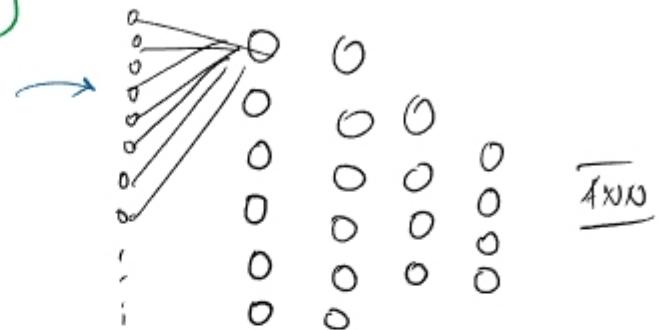
CNN Vs ANN

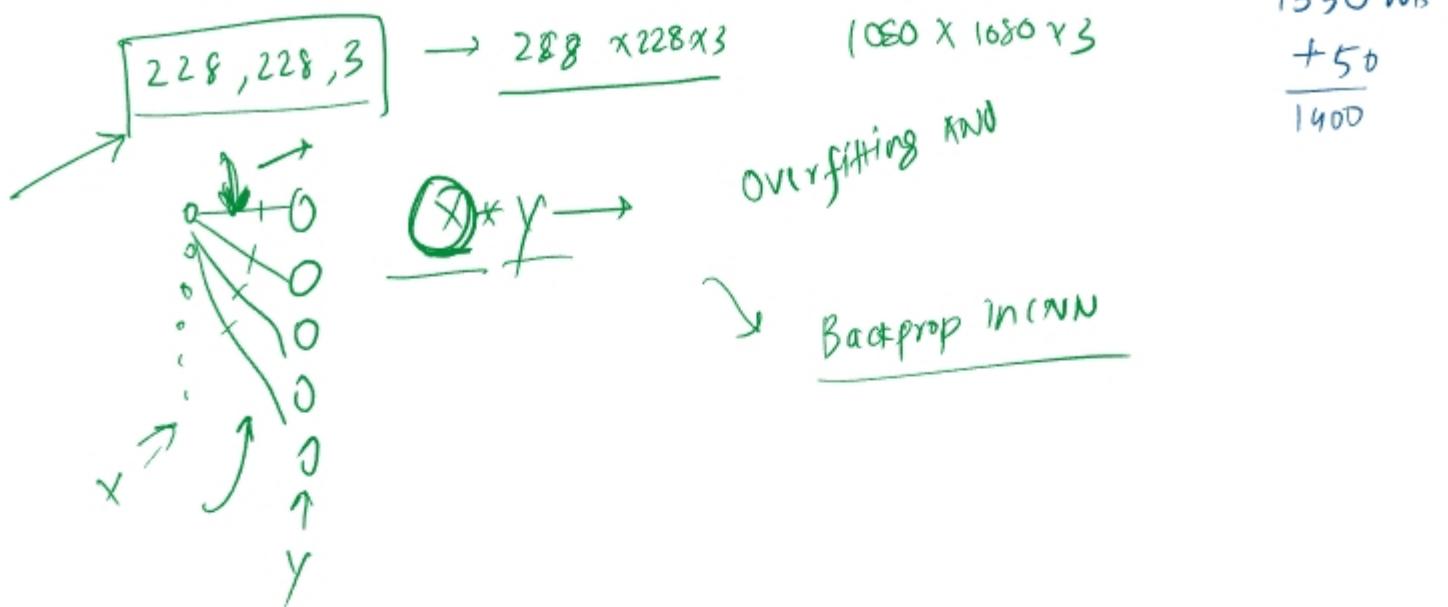
06 September 2022 10:00

- 1) Computation Cost → W
- 2) Overfitting →
- 3) Loss of imp features like spatial arrangement of pixels



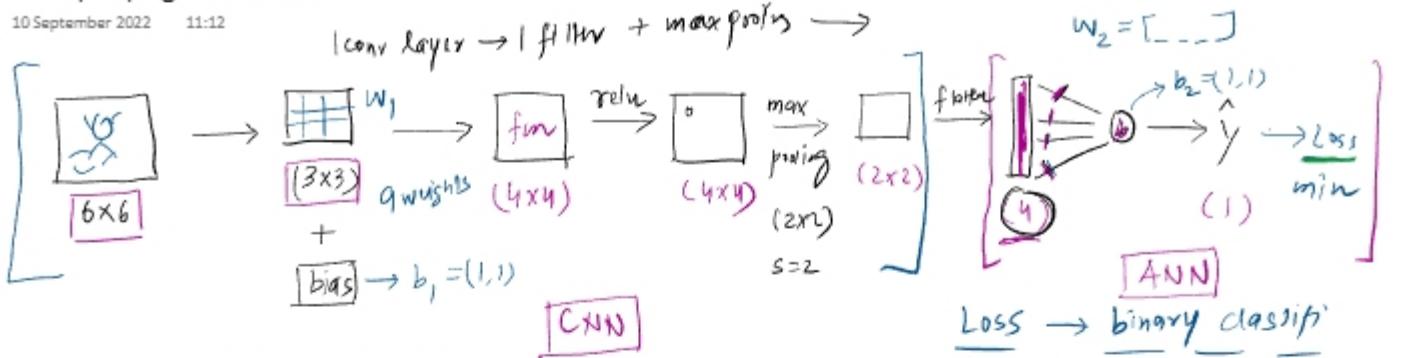
2D →





Backpropagation in CNN

10 September 2022 11:12



Trainable Parameters

$$\begin{aligned} W_1 &= (3, 3) & W_2 &= (1, 4) \\ b_1 &= (1, 1) & b_2 &= (1, 1) \end{aligned} = 15 \text{ trainable parameters}$$

Logical Flow

Forward Prop

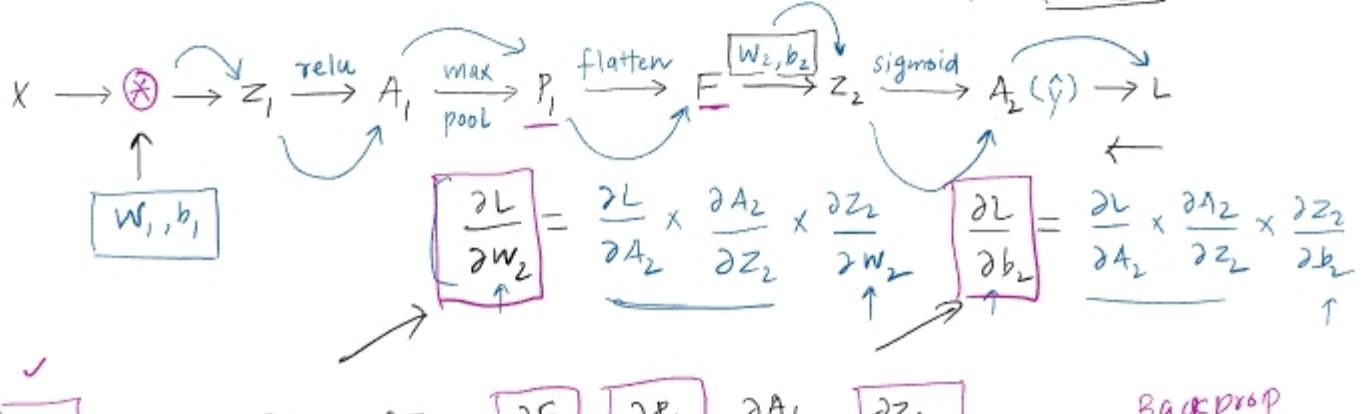
$$\left\{ \begin{array}{l} z_1 = \text{conv}(x, W_1) + b_1 \\ A_1 = \text{relu}(z_1) \\ P_1 = \text{maxpool}(A_1) \\ F = \text{flatten}(P_1) \\ z_2 = W_2 F + b_2 \\ A_2 = \sigma(z_2) \end{array} \right.$$

Gradient Descent

$$W_1 = W_1 - \eta \frac{\partial L}{\partial W_1} \quad W_2 = W_2 - \eta \frac{\partial L}{\partial W_2} \quad \text{Loss is minimized}$$

$$b_1 = b_1 - \eta \frac{\partial L}{\partial b_1}$$

$$b_2 = b_2 - \eta \frac{\partial L}{\partial b_2}$$

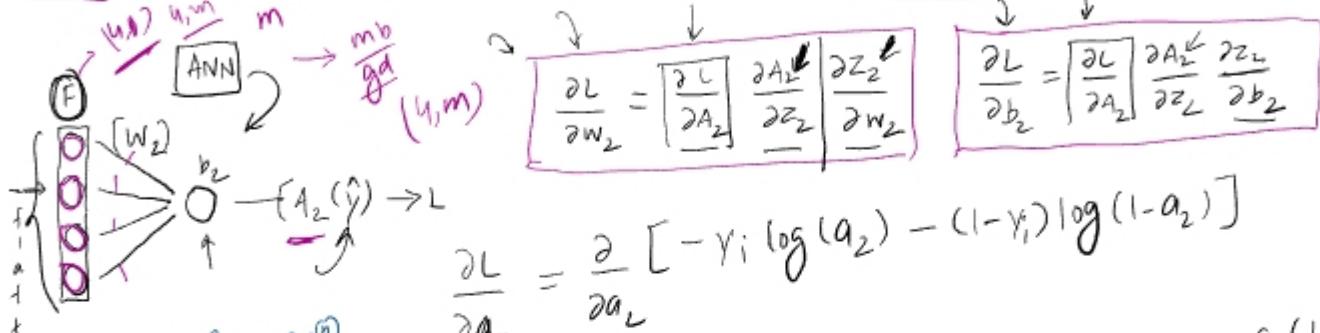


$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial F} \times \boxed{\frac{\partial F}{\partial P_1}} \times \boxed{\frac{\partial P_1}{\partial A_1}} \times \frac{\partial A_1}{\partial Z_1} \times \boxed{\frac{\partial Z_1}{\partial w_1}}$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial F} \times \frac{\partial F}{\partial P_1} \times \frac{\partial P_1}{\partial A_1} \times \frac{\partial A_1}{\partial Z_1} \times \boxed{\frac{\partial Z_1}{\partial b_1}}$$

Backprop

→ Convolution
→ Flatten
→ Max pooling



$$\frac{\partial L}{\partial w_2} = \boxed{\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial Z_2} \frac{\partial Z_2}{\partial w_2}}$$

$$\frac{\partial L}{\partial b_2} = \boxed{\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial Z_2} \frac{\partial Z_2}{\partial b_2}}$$

$$\frac{\partial L}{\partial a_2} = \frac{\partial}{\partial a_2} [-y_i \log(a_2) - (1-y_i) \log(1-a_2)]$$

$$\frac{\partial L}{\partial a_2}$$

$$(1, m) = -\frac{y_i}{a_2} + \frac{(1-y_i)}{(1-a_2)} = \frac{-y_i(1-a_2) + a_2(1-y_i)}{a_2(1-a_2)}$$

$$\frac{\partial L}{\partial a_2} = -y_i + \frac{a_2 - y_i}{a_2(1-a_2)} = \frac{(a_2 - y_i)}{a_2(1-a_2)}$$

$$\frac{\partial A_2}{\partial Z_2} = \sigma(z_2) [1 - \sigma(z_2)] = \underline{a_2[1-a_2]}$$

$$\boxed{\frac{\partial Z_2}{\partial w_2} = F}$$

$$\frac{\partial Z_2}{\partial b_2} = 1$$

w₂ update
shape =

$$\frac{\partial L}{\partial w_2} = \frac{(a_2 - y_i)}{a_2(1-a_2)} \times a_2(1-a_2) \times F = (a_2 - y_i) F = \frac{(A_2 - Y) F^T}{(1,1) \quad (1,1)}$$

$$\frac{\partial L}{\partial b_2} = \frac{(a_2 - y_i)}{a_2(1-a_2)} \times a_2(1-a_2) \times 1 = (A_2 - Y)$$

m images

(1,1) (1,9)
(1,9)

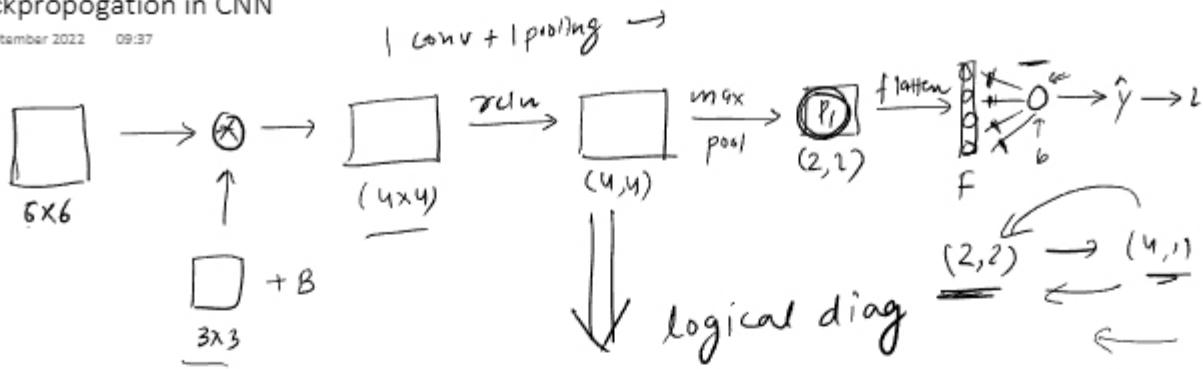
$$\boxed{\frac{\partial L}{\partial b_2} = (A_2 - Y) F^T} \quad \boxed{\frac{\partial L}{\partial b_2} = (A_2 - Y)}$$

$$\left[\frac{\partial L}{\partial w_2} \right] = (A_2 - Y)^\top \quad \left[\frac{\partial L}{\partial b_2} \right] = -1$$

$\uparrow \quad (1, m) - (1, m)$
 $(1, m) \quad (m, n) \rightarrow \boxed{(1, 4)} \rightarrow \begin{array}{l} \text{image} \\ \text{batch of image} \end{array}$
 $w_2 \rightarrow (1, n)$

Backpropagation in CNN

15 September 2022 09:37



Forward Prop

$$Z_1 = \text{conv}(X, w_1) + b_1$$

$$A_1 = \text{relu}(Z_1)$$

$$P_1 = \text{maxpool}(A_1)$$

$$F = \text{flatten}(P_1)$$

$$z_2 = w_2 F + b_2$$

$$A_2 = \sigma(z_2)$$

$$L = \frac{1}{m} \sum_{i=1}^m [-y_i \log(A_2) - (1-y_i) \log(1-A_2)]$$

6 derivatives

$$\text{conv}(X, \frac{\partial L}{\partial z_1})$$

$$\left[\frac{\partial z_2}{\partial F} \right] = W_2 \rightarrow$$

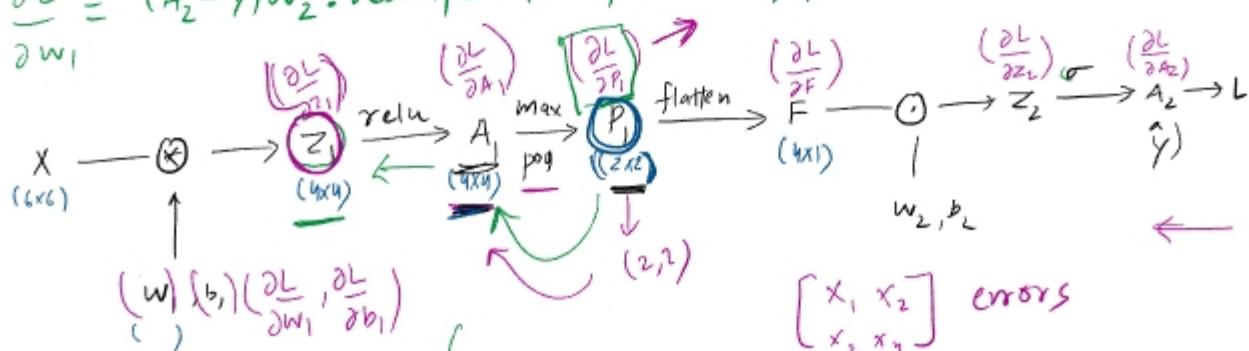
$$\text{Shape?} \rightarrow (F)$$

$$\frac{\partial F}{\partial P_1} \quad \text{no trainable parameters}$$

reshape(P1.shape)

$$\frac{\partial L}{\partial w_1} = (A_2 - y) W_2 \cdot \text{reshape}(P_1, \text{shape})$$

$$\frac{\partial L}{\partial b_1}$$

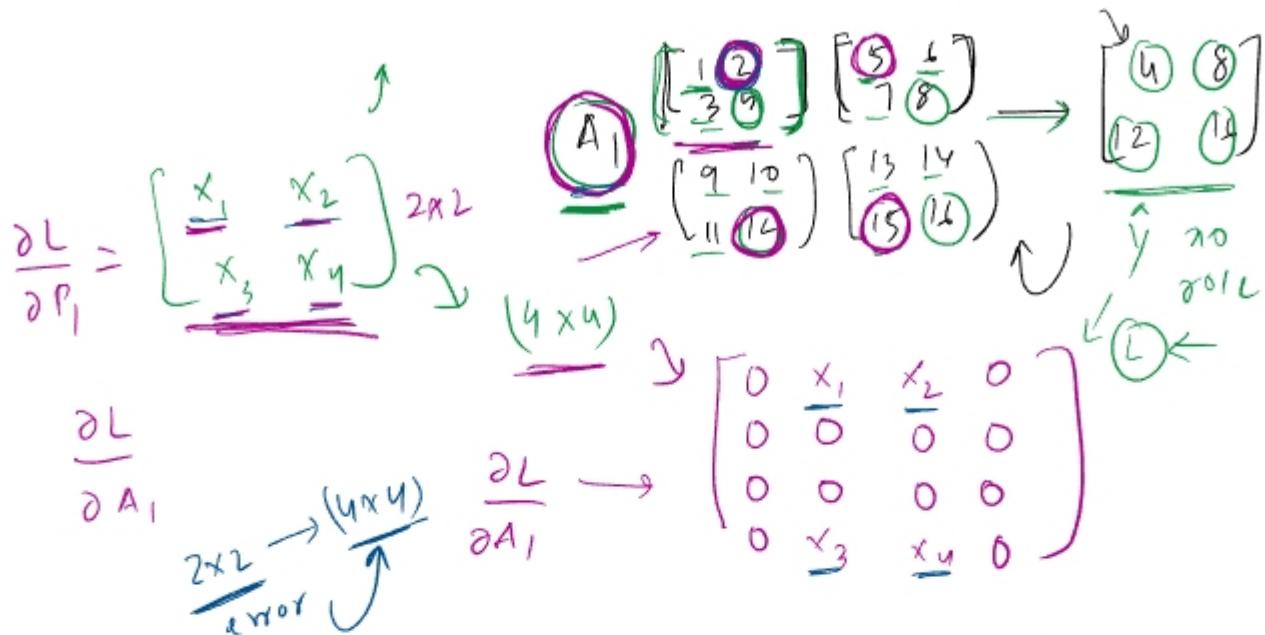


$$\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix} \text{ errors}$$

$$\frac{\partial L}{\partial A_1} = (4, 4)$$

$$\frac{\partial L}{\partial A_1} = (4, 4)$$

flatten \rightarrow no trainable parameters



A, P

$$\frac{\partial L}{\partial W_L} = \left[\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial Z_2} \frac{\partial Z_2}{\partial F} \frac{\partial F}{\partial P_1} \frac{\partial P_1}{\partial A_1} \frac{\partial A_1}{\partial Z_1} \frac{\partial Z_1}{\partial W_1} \right]$$

$(A_2 - y) W_2 \cdot \text{reshape}(P_1, \text{shape})$

$$\frac{\partial L}{\partial b_1} = \left[\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial Z_2} \frac{\partial Z_2}{\partial F} \frac{\partial F}{\partial P_1} \frac{\partial P_1}{\partial A_1} \frac{\partial A_1}{\partial Z_1} \frac{\partial Z_1}{\partial b_1} \right] \rightarrow \frac{\partial L}{\partial A_1}$$

$$\frac{\partial L}{\partial P_1} \quad \frac{\partial L}{\partial A_1}$$

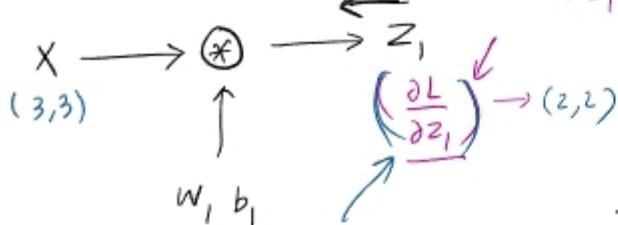
$\frac{\partial L}{\partial A_1} = \begin{cases} \frac{\partial L}{\partial P_1}_{xy}, & \text{if } A_{mn} \text{ is the max element} \\ 0, & \text{otherwise} \end{cases}$

$\frac{\partial A_1}{\partial Z_1} = \begin{cases} 1 & \text{if } Z_{1xy} > 0 \\ 0 & \text{if } Z_{1xy} < 0 \end{cases}$

$\left(\frac{\partial L}{\partial Z_1} \right)$

Convolution \hookrightarrow Backprop \hookrightarrow max pooling \hookrightarrow flatten

Backprop on Convolution



$$\frac{\partial L}{\partial Z_1} = \begin{bmatrix} \frac{\partial L}{\partial Z_{11}} & \frac{\partial L}{\partial Z_{12}} \\ \frac{\partial L}{\partial Z_{21}} & \frac{\partial L}{\partial Z_{22}} \end{bmatrix}$$



$$X = \underbrace{\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{bmatrix}}_{\text{---}} \times_{\textcircled{*}} \underbrace{\begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}}_{\text{---}} + b_1$$

$$\begin{bmatrix} \frac{\partial L}{\partial b_1} \end{bmatrix} = \frac{\partial L}{\partial z_1} \times \frac{\partial z_1}{\partial b_1} = \left[\frac{\frac{\partial L}{\partial z_{11}} \times \frac{\partial z_{11}}{\partial b_1}}{\uparrow} + \frac{\frac{\partial L}{\partial z_{12}} \times \frac{\partial z_{12}}{\partial b_1}}{\uparrow} + \frac{\frac{\partial L}{\partial z_{21}} \times \frac{\partial z_{21}}{\partial b_1}}{\uparrow} + \frac{\frac{\partial L}{\partial z_{22}} \times \frac{\partial z_{22}}{\partial b_1}}{\uparrow} \right]$$

$$= \left(\frac{\partial L}{\partial z_{11}} + \frac{\partial L}{\partial z_{12}} + \frac{\partial L}{\partial z_{21}} + \frac{\partial L}{\partial z_{22}} \right) = \text{sum} \left(\frac{\partial L}{\partial z_i} \right)$$

$$\begin{bmatrix} \frac{\partial L}{\partial b_1} \end{bmatrix} = \text{sum} \left(\frac{\partial L}{\partial z_i} \right) \rightarrow \text{scalar}$$

bias

$$X \rightarrow \textcircled{*} \rightarrow \textcircled{z}_1$$

\uparrow

(3×3)

\uparrow

w_1, b_1

(2×2)

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \quad w_1 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} + \text{?}$$

$$\frac{\partial L}{\partial w_1} = \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} \end{bmatrix} \quad \frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z_{11}} \times \boxed{\frac{\partial z_{11}}{\partial w_{11}}} + \frac{\partial L}{\partial z_{12}} \times \boxed{\frac{\partial z_{12}}{\partial w_{11}}} + \frac{\partial L}{\partial z_{21}} \times \boxed{\frac{\partial z_{21}}{\partial w_{11}}} + \frac{\partial L}{\partial z_{22}} \times \boxed{\frac{\partial z_{22}}{\partial w_{11}}}$$

$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial z_{11}} \times \boxed{\frac{\partial z_{11}}{\partial w_{12}}} + \frac{\partial L}{\partial z_{12}} \times \boxed{\frac{\partial z_{12}}{\partial w_{12}}} + \frac{\partial L}{\partial z_{21}} \times \boxed{\frac{\partial z_{21}}{\partial w_{12}}} + \frac{\partial L}{\partial z_{22}} \times \boxed{\frac{\partial z_{22}}{\partial w_{12}}}$$

$$\frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial z_{11}} \underbrace{\frac{\partial z_{11}}{\partial w_{21}}} + \frac{\partial L}{\partial z_{12}} \underbrace{\frac{\partial z_{12}}{\partial w_{21}}} + \frac{\partial L}{\partial z_{21}} \underbrace{\frac{\partial z_{21}}{\partial w_{21}}} + \frac{\partial L}{\partial z_{22}} \underbrace{\frac{\partial z_{22}}{\partial w_{21}}}$$

$$\frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial z_{11}} \underbrace{\frac{\partial z_{11}}{\partial w_{22}}} + \frac{\partial L}{\partial z_{12}} \underbrace{\frac{\partial z_{12}}{\partial w_{22}}} + \frac{\partial L}{\partial z_{21}} \underbrace{\frac{\partial z_{21}}{\partial w_{22}}} + \frac{\partial L}{\partial z_{22}} \underbrace{\frac{\partial z_{22}}{\partial w_{22}}}$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial z_{11}} x_{11} + \frac{\partial L}{\partial z_{12}} x_{12} + \frac{\partial L}{\partial z_{21}} x_{21} + \frac{\partial L}{\partial z_{22}} x_{22} \\ \frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial z_{11}} x_{12} + \frac{\partial L}{\partial z_{12}} x_{13} + \frac{\partial L}{\partial z_{21}} x_{22} + \frac{\partial L}{\partial z_{22}} x_{23} \\ \frac{\partial L}{\partial w_{21}} = \frac{\partial L}{\partial z_{11}} x_{21} + \frac{\partial L}{\partial z_{12}} x_{22} + \frac{\partial L}{\partial z_{21}} x_{31} + \frac{\partial L}{\partial z_{22}} x_{32} \\ \frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial z_{11}} x_{22} + \frac{\partial L}{\partial z_{12}} x_{23} + \frac{\partial L}{\partial z_{21}} x_{32} + \frac{\partial L}{\partial z_{22}} x_{33} \end{array} \right. \quad \left. \begin{array}{l} X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \quad \frac{\partial L}{\partial z_i} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \\ \frac{\partial L}{\partial z_{31}} & \frac{\partial L}{\partial z_{32}} \end{bmatrix} \\ \frac{\partial L}{\partial w_i} = \text{conv}(X, \frac{\partial L}{\partial z_i}) \end{array} \right.$$

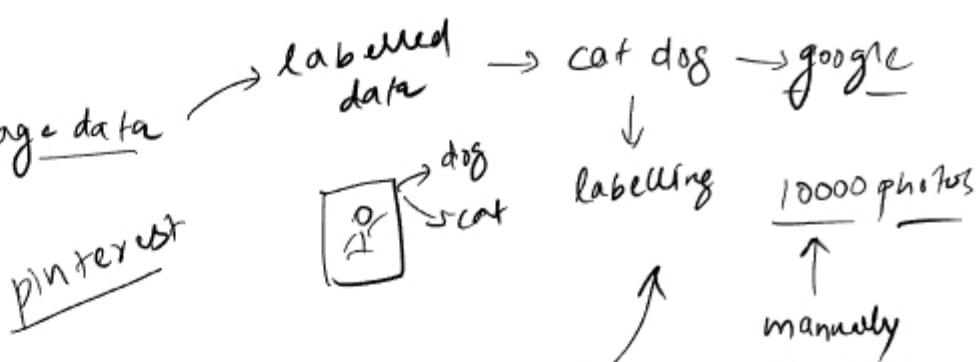
$$\frac{\partial L}{\partial w_1} = \text{conv}(X, \frac{\partial L}{\partial z_1})$$

$$\frac{\partial L}{\partial b_1} = \text{sum}(\frac{\partial L}{\partial z_1})$$

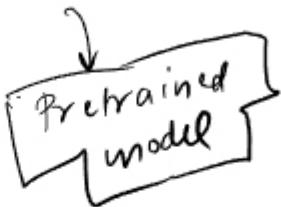
Why use Pretrained models?

03 October 2022 12:52

1) Data hungry → image data



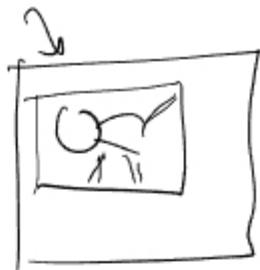
2) Time → model → training



ImageNET Dataset

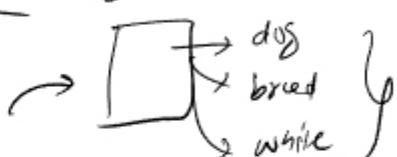
03 October 2022 12:36

Visual Database of images (Why What and How)



Why

2006 → Fei fei Li → model and algorithms

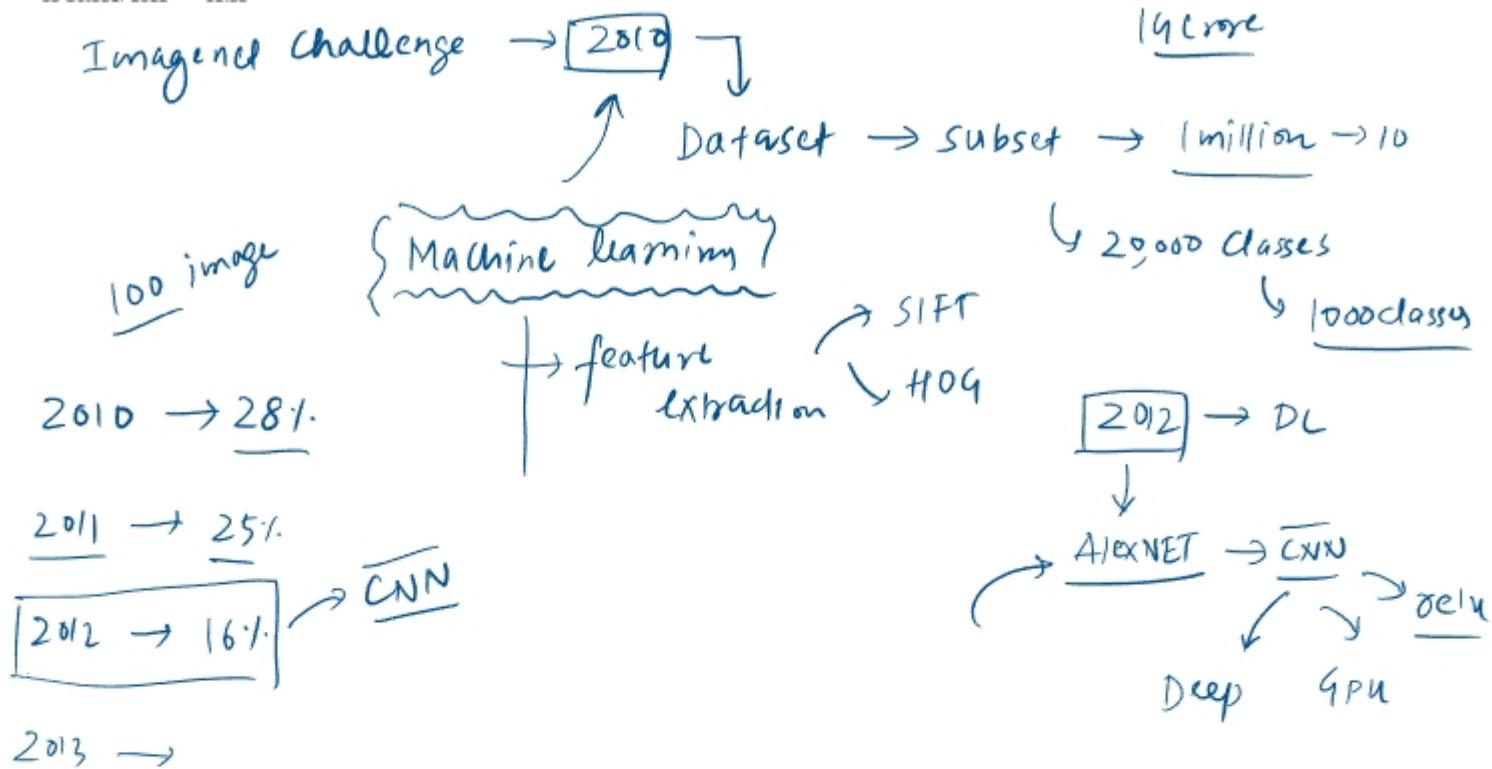
1.4 million images → 1.4 crore image → 20,000 categories → labels
1 million images → to batch → bounding box
↓ object local → 

How → crowd sourcing

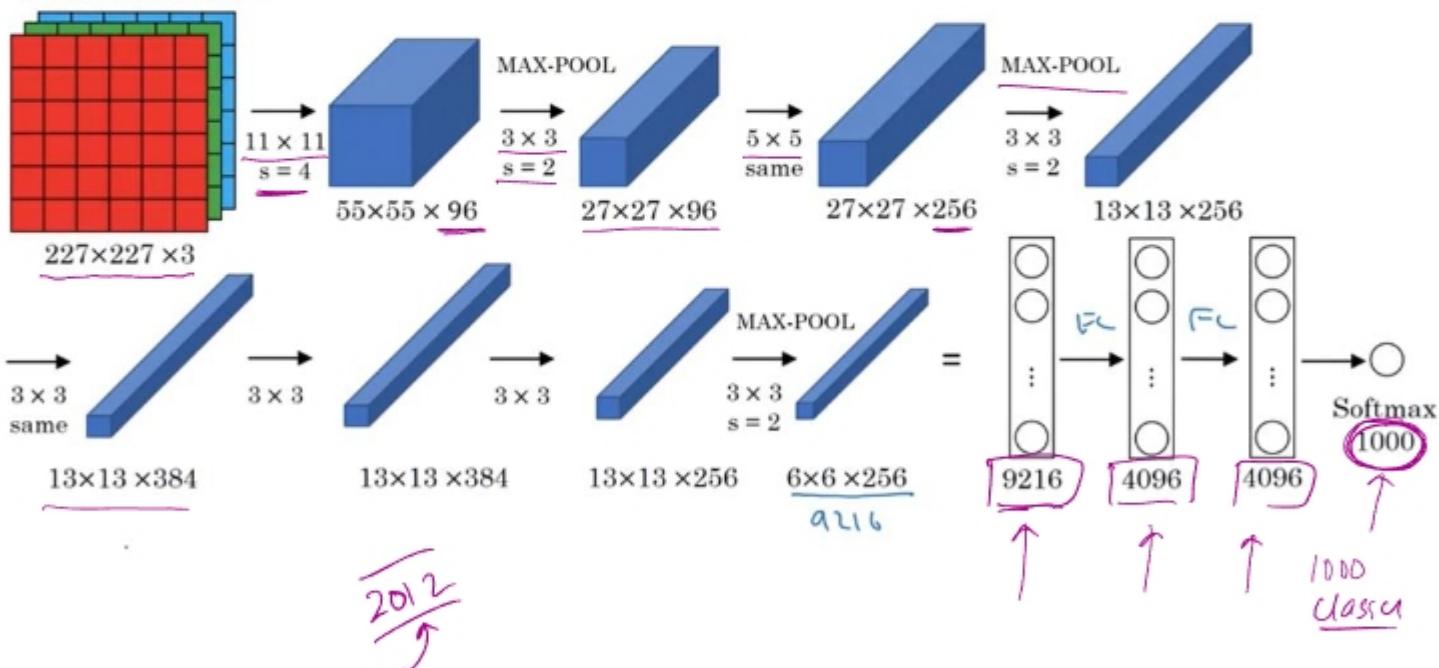
 → Amazon mechanical Turk

→ Datasup → Deep learning

ImageNET challenge



AlexNet

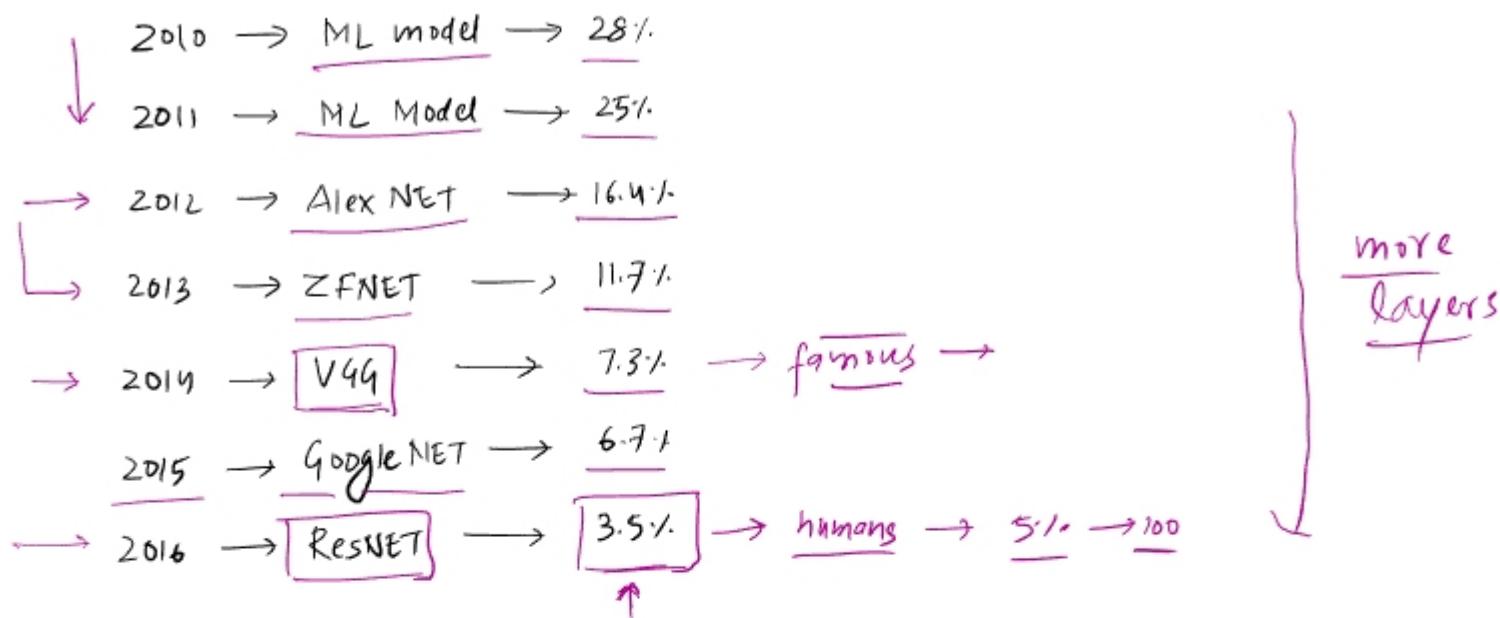


[Krizhevsky et al., 2012. ImageNet classification with deep convolutional neural networks]

Andrew Ng

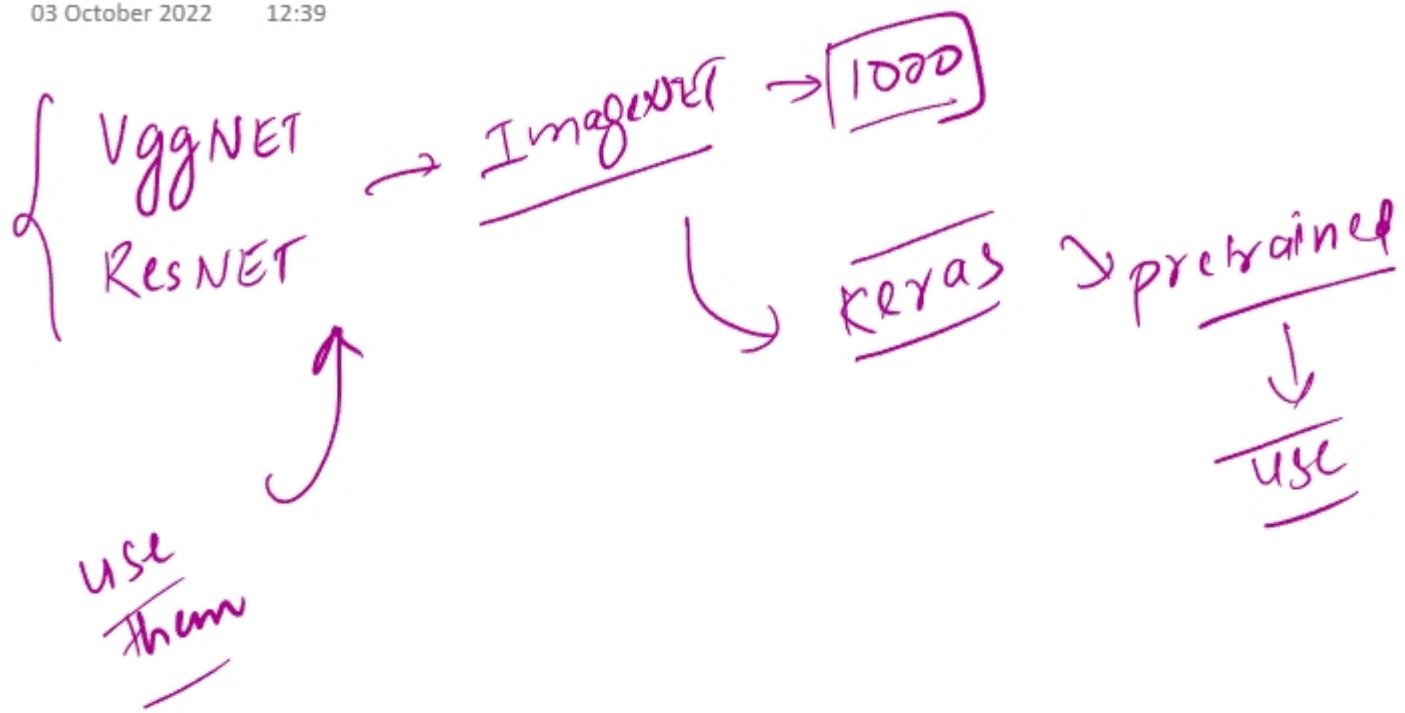
Famous Architectures

03 October 2022 12:39



Idea of Pretrained Models

03 October 2022 12:39



Keras Demo

03 October 2022 12:40

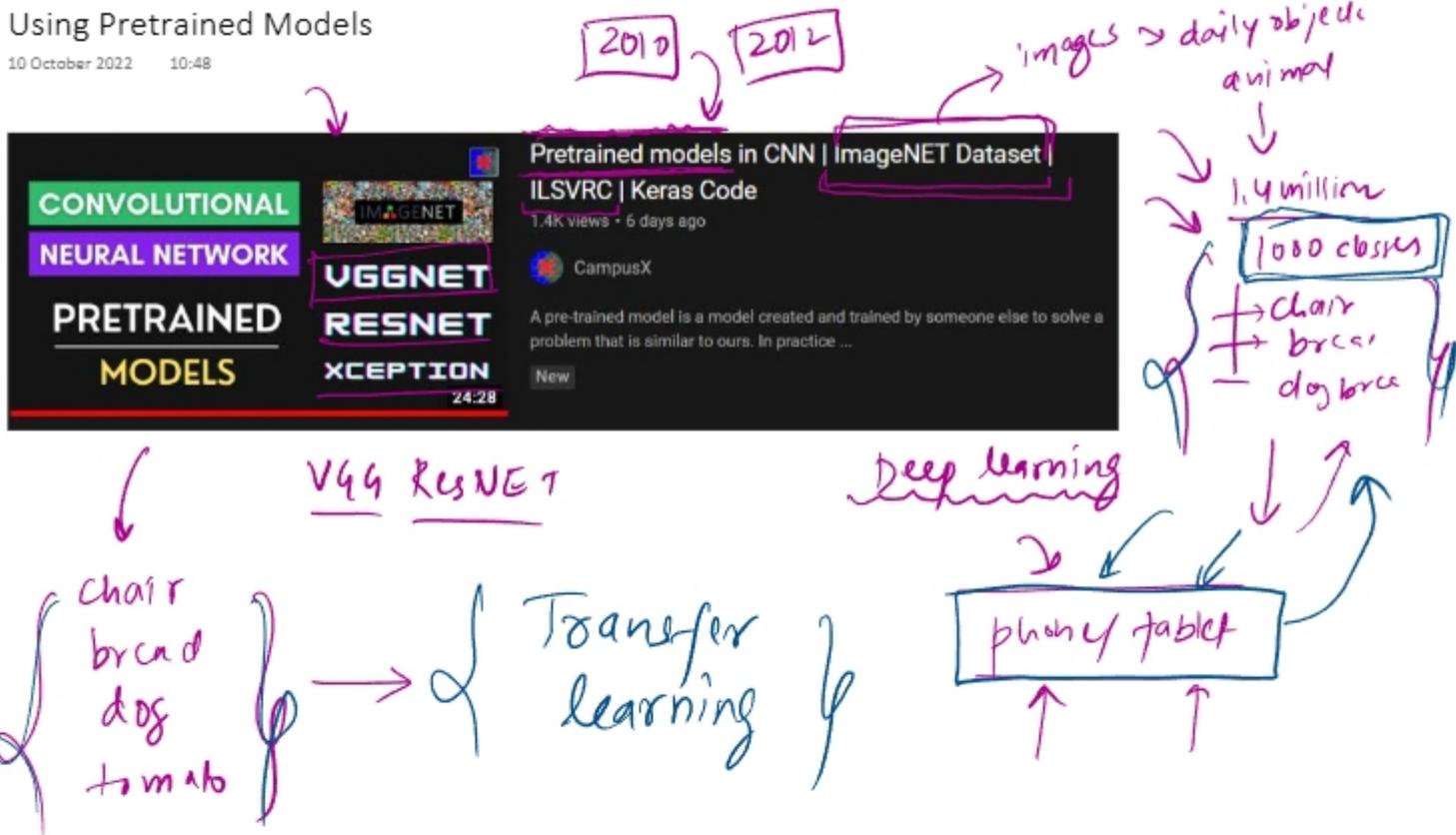
Problem with training your own model

10 October 2022 10:48

- 1) Data hungry → labelled → 10,000 → google
↓
cat/dog → manual labor
- 2) lot of time → days/weeks

Using Pretrained Models

10 October 2022 10:48



Transfer Learning

10 October 2022 10:49

Transfer learning is a research problem in machine learning that focuses on storing knowledge gained while solving one problem and applying it to a different but related problem.

