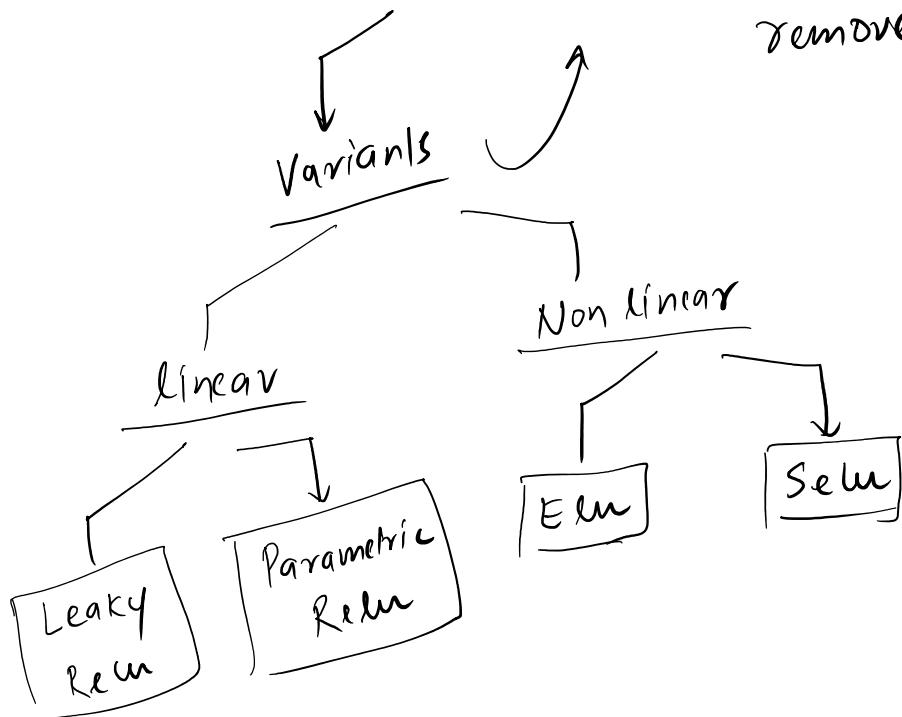


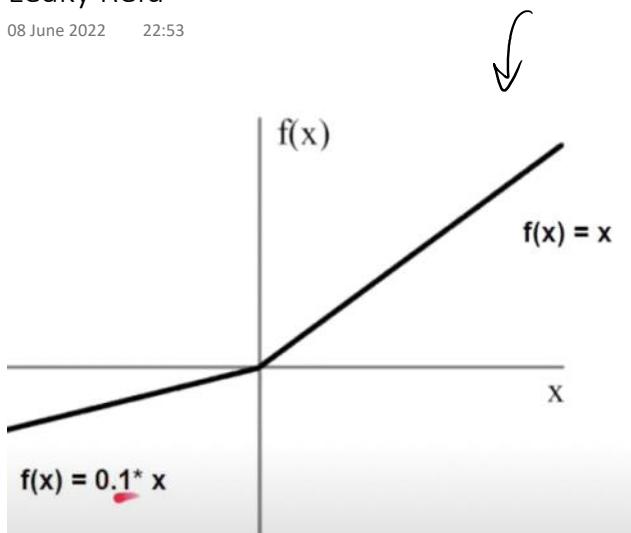
Solutions

- Set low learning rate
- bias $\rightarrow +ve$ value $\rightarrow 0.01$
- Don't use ReLU \rightarrow variants (+ve keep)
remove



Leaky Relu

08 June 2022 22:53



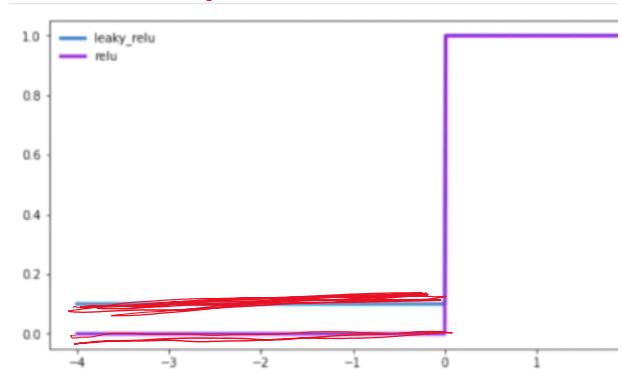
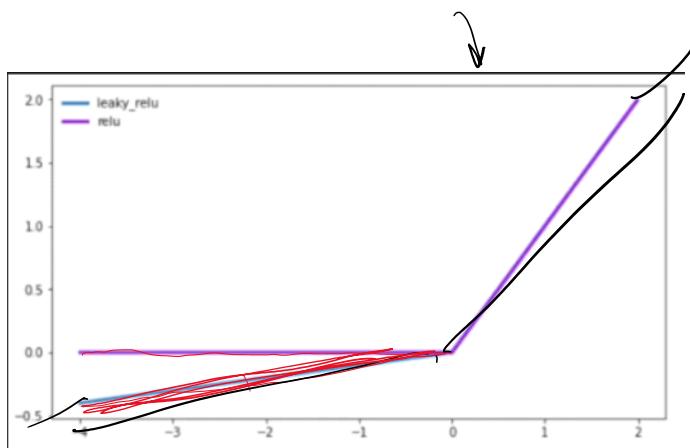
$$f'(z) = \max(0.01 z, z)$$

$\frac{\partial}{\partial z} f(z) = \begin{cases} 1 & z \geq 0 \\ 0.01 & z < 0 \end{cases}$

$-ve(z) \rightarrow \neq 0 \quad \frac{\partial}{\partial z} z$

\hookrightarrow derivative

$$f'(z) = \begin{cases} 1 & z \geq 0 \\ 0.01 & z < 0 \end{cases}$$



Advantages

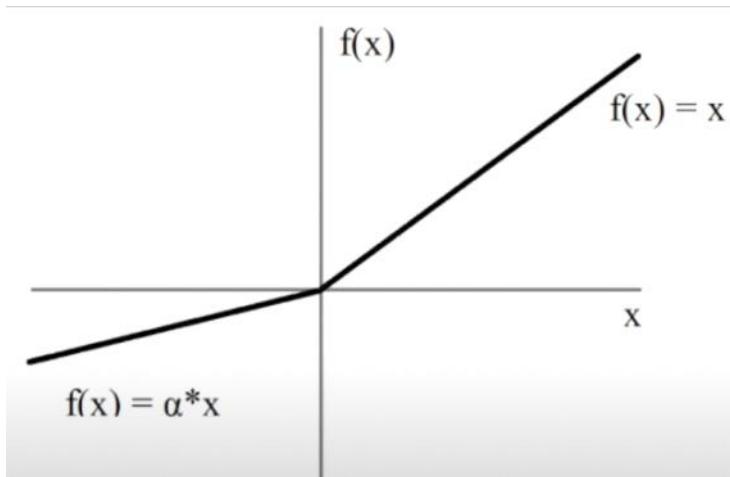
- Non-saturated \rightarrow Unbounded
- Easily computed
- No dying relu problem
- Close to 0 centered

Disadv

-ve/+ve

Parametric Relu

08 June 2022 22:53



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \underline{\alpha x} & \text{otherwise} \\ 0.0x \end{cases}$$

$\alpha \rightarrow$ trainable parameter

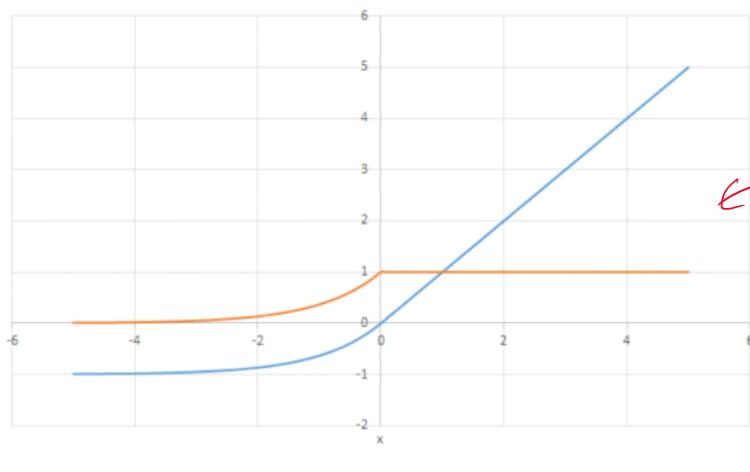
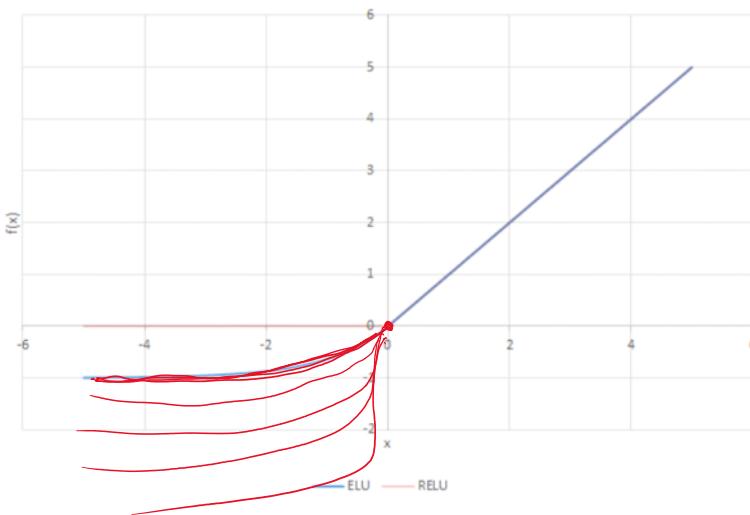
flexibility

depend on data

Elu - Exponential Linear Unit

09 June 2022 00:29

performance better Relu



$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

0.1 to 0.3

$$\text{ELU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ \text{ELU}(x) + \alpha & \text{if } x \leq 0 \end{cases}$$

Advantages

→ Close to zero centered
↳ convergence faster

- Better generalized
- Dying ReLU x
- Always continuous as well as differentiable

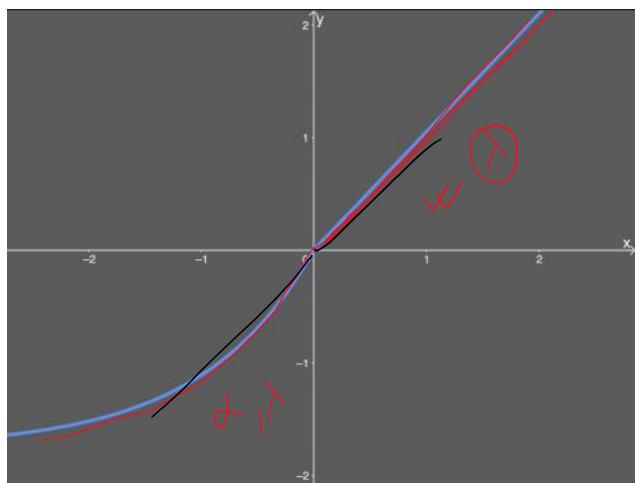
Disadv

→ computation expensive $\circledcirc @x$

Selu - Scaled Exponential Linear Unit

09 June 2022 00:29

Recent \rightarrow (new) \rightarrow 90+ pages



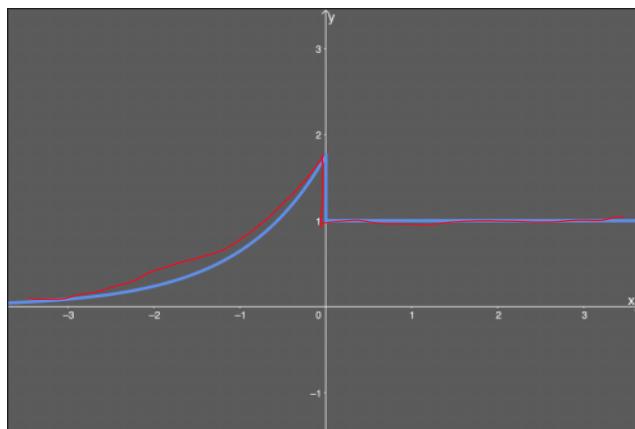
$$\text{SELU}(x) = \lambda \begin{cases} x & \text{if } x > 0 \\ \alpha e^x - \alpha & \text{if } x \leq 0 \end{cases}$$

$$a \approx 1.6732632423543772848170429916717$$

$$\lambda \approx 1.0507009873554804934193349852946$$

\rightarrow fixed
train bnx

$$\text{SELU}'(x) = \lambda \begin{cases} 1 & \text{if } x > 0 \\ \alpha e^x & \text{if } x \leq 0 \end{cases}$$



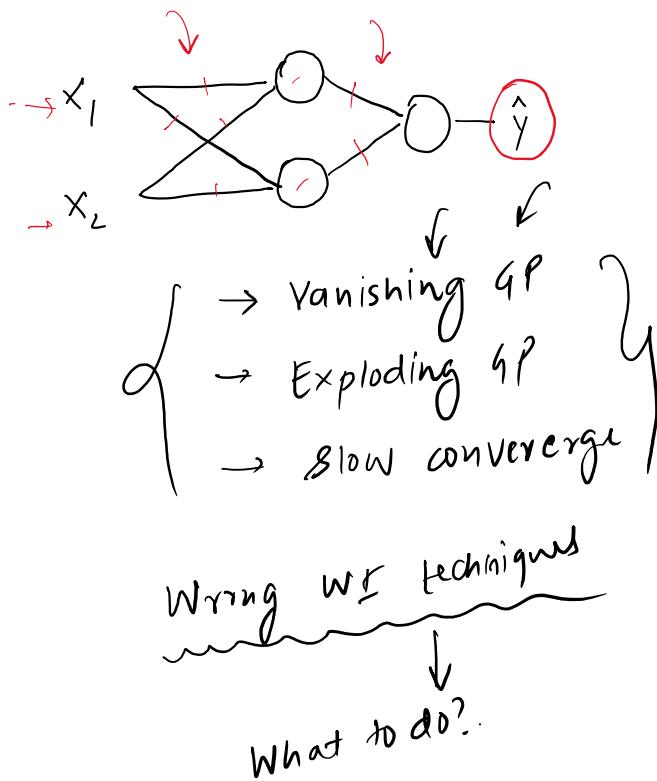
Self-normalizing \rightarrow activation
 \downarrow
normalized

$m = 0 \quad \sigma = 1$

converge faster

Why Weight Initialization is Important?

22 June 2022 12:48



Deep
Learning
2000's

1. Initialize the parameters
2. Choose an *optimization algorithm* ✓
3. Repeat these steps:

1. Forward propagate an input
2. Compute the cost function
3. Compute the gradients of the cost with respect to parameters using backpropagation
4. Update each parameter using the gradients, according to the optimization algorithm

Sigmoid
Activ

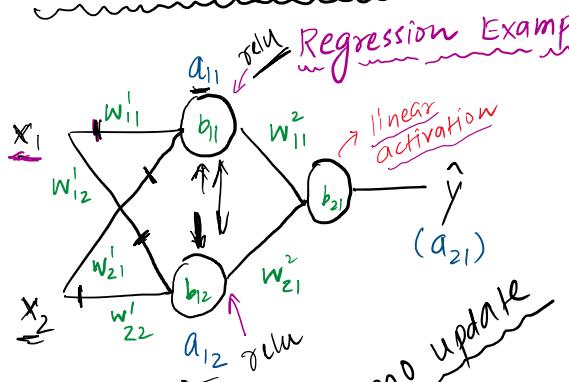
Wrong
weigh
init

What to do?

What not to do?

22 June 2022 12:49

Case 1 → zero Initialization



$$w_{11}^1 = \underline{w_{11}^1} - \eta \frac{\partial L}{\partial w_{11}^1}$$

→ relu
→ tanh
→ sigmoid

$$w=0, b=0$$

$$a_{11} = \frac{e^{z_{11}} - e^{-z_{11}}}{e^{z_{11}} + e^{-z_{11}}} = 0$$

$$a_{11} = \max(0, z_{11}) = 0$$

$$z_{11} = \underline{w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}} = 0$$

$$a_{12} = \max(0, z_{12}) = 0$$

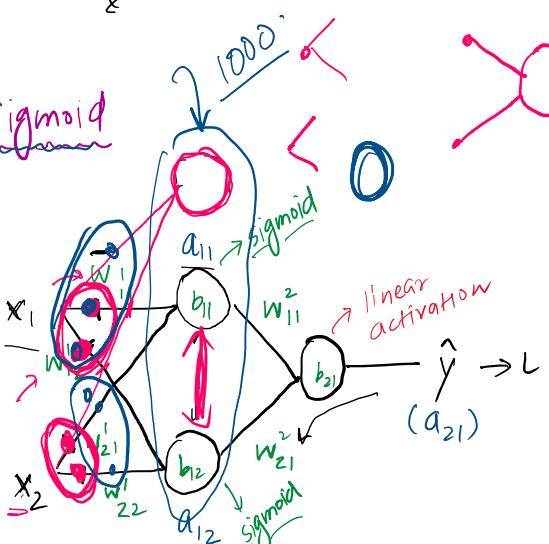
$$z_{12} = \underline{w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}} = 0$$

$$\boxed{a_{11} = a_{12}} \rightarrow \text{equal and } 0$$

$$w=0$$

No training will place

$$f$$



$$a_{11} = \sigma(z_{11}) = 0.5$$

$$z_{11} =$$

$$a_{12} = 0.5$$

$$a_{11} = a_{12}$$

Non-linear part

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

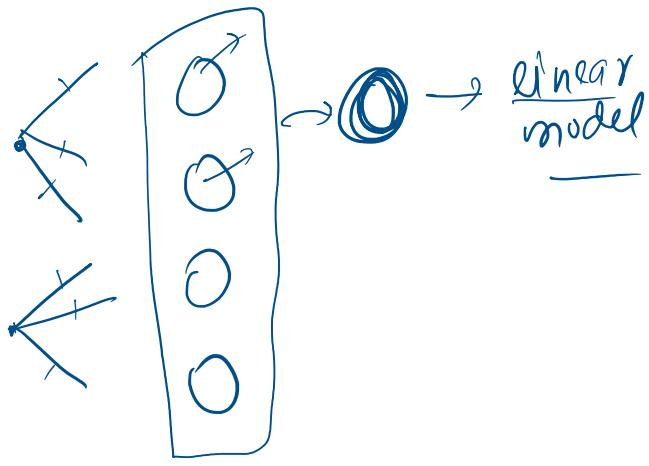
$$a_{11} = a_{12}, z_{11} = z_{12}$$

$$\frac{\partial L}{\partial w_{11}^1} = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}} & \frac{\partial \hat{y}}{\partial a_{11}} & \frac{\partial a_{11}}{\partial z_{11}} \end{bmatrix} X_1$$

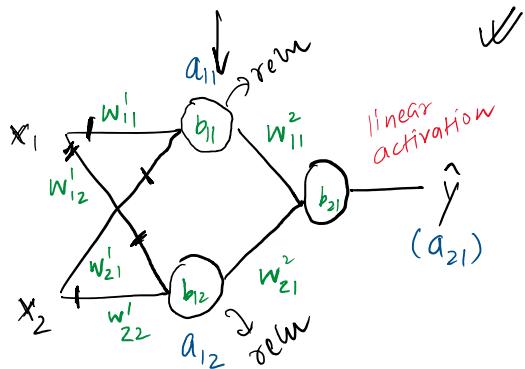
$$\frac{\partial L}{\partial w_{12}^1} = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}} & \frac{\partial \hat{y}}{\partial a_{12}} & \frac{\partial a_{12}}{\partial z_{12}} \end{bmatrix} X_1$$

$$\frac{\partial L}{\partial w_{21}^1} = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}} & \frac{\partial \hat{y}}{\partial a_{11}} & \frac{\partial a_{11}}{\partial z_{11}} \end{bmatrix} X_2$$

$$\frac{\partial L}{\partial w_{22}^1} = \begin{bmatrix} \frac{\partial L}{\partial \hat{y}} & \frac{\partial \hat{y}}{\partial a_{12}} & \frac{\partial a_{12}}{\partial z_{12}} \end{bmatrix} X_2$$



Case 2 → Non-0 constant value



$$w = 0.5 \quad b = 0.5 \quad \text{sigmoid} \rightarrow 0$$

$$a_{11} = \max(0, z_{11}) \rightarrow \text{non-zero}$$

$$z_{11} = w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11} \neq 0$$

$$a_{12} = \max(0, z_{12}) \neq 0$$

$$z_{12} = w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12} \neq 0$$

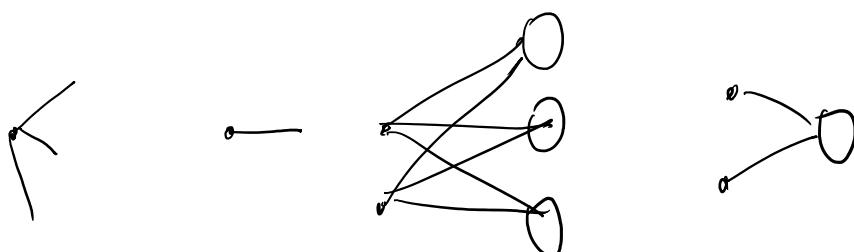
$$z_{12} = z_{11} \quad \boxed{a_{11} = a_{12}}$$

$$\boxed{\frac{\partial L}{\partial w_{11}^1}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} x_1$$

$$\boxed{\frac{\partial L}{\partial w_{12}^1}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} x_1$$

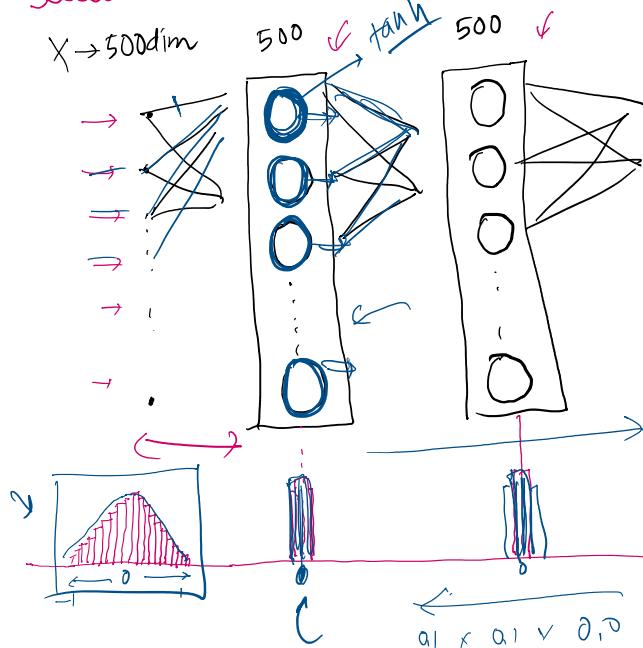
$$\frac{\partial L}{\partial w_{21}^1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{11}} \frac{\partial a_{11}}{\partial z_{11}} x_2$$

$$\frac{\partial L}{\partial w_{22}^1} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{12}} \frac{\partial a_{12}}{\partial z_{12}} x_2$$



of random init ↗

Case 3 - Random Initialization



VGP sigmoid

$x_1, \dots, x_{500} | y$

gaussian distribution

weight init → random
→ small values

$\text{np.random.randn}(500, 500) \times 0.01$

$\boxed{\text{bias} = 0}$

$x_i \rightarrow [0-1] \rightarrow [-1+0]$

$\sum w_i x_i$

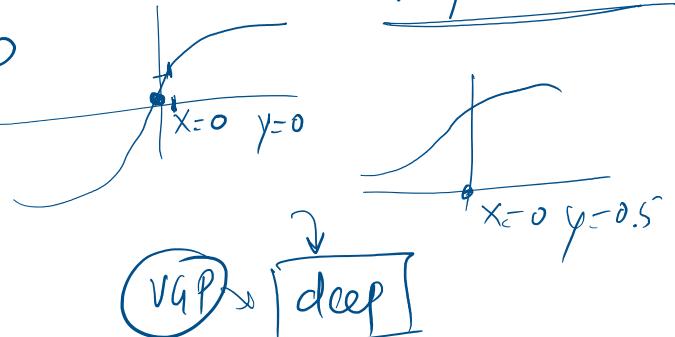
$[0.0007]$

small number tanh

② is a small
very close to 0

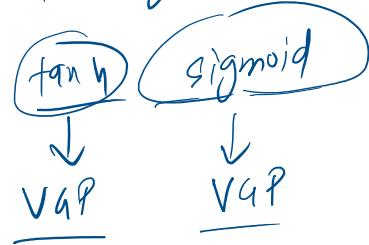
$w = w - \eta \frac{\partial L}{\partial w} \approx 0$

② vanishing gradients problem

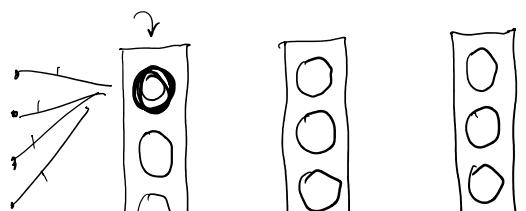


ReLU ↗
small values
→ slow convergence

small weights → 0.0007



Random Initialization (Large values)

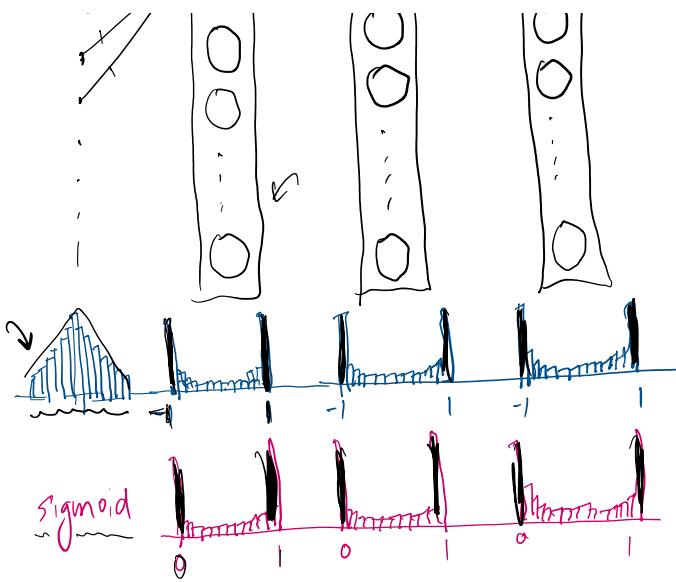


$\text{np.random.randn}(500, 500) \rightarrow [0-1]$

$\sum w_i x_i \rightarrow 500 \rightarrow [100 \rightarrow 500]$

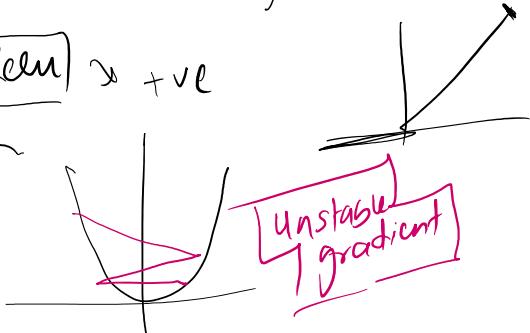
$0-1$

tanh/sigmoid



• 0-1 tanh / sigmoid
 $\sum w_i x_i \rightarrow g(z)$
 saturation $\rightarrow 100$ \rightarrow saturate
 $\Sigma w_i x_i \rightarrow$ $\frac{p_{new}}{250} \downarrow \frac{1}{250}$
 slow training
 vanishing gradient problem
 gradients

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

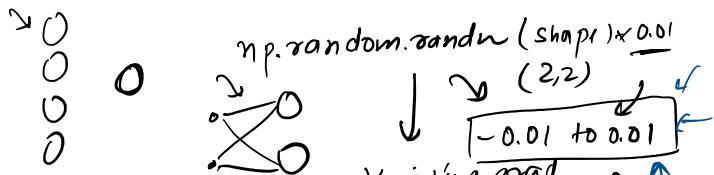


- zero x }
- non-zero same x }
- small random /
- large random x }

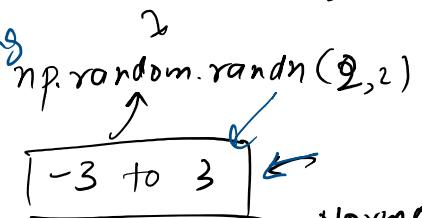
Key Insights

22 June 2022 12:49

- Things not to do -
- 1) Zero initialization → no training
- 2) Non-zero constant initialization → linear
- 3) Random initialization with small weights → vanishing/exploding gradients
- 4) Random initialization with large weights → random

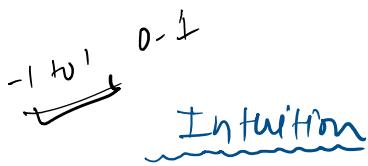


Vanishing grad
Slow convergence



What can be done

→ Heuristics (Jugad) → practical soln



-0.1 to 0.1

small ←

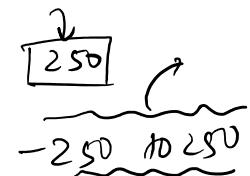
np.random.randn(250, 250) * 0.01

large ←

np.random.randn(250, 250) * 1 → -3 to 3

$$\text{Variance} = \frac{1}{n} = \frac{1}{250}$$

$$* \sqrt{\frac{1}{250}} \rightarrow [sd]$$



∴

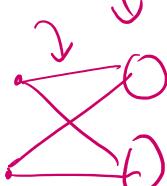
$$n \rightarrow 1 \downarrow \boxed{\frac{1}{n}} \downarrow \uparrow \sigma_w$$

Xavier init
(Normal)

$$\sqrt{\frac{1}{\text{fan-in}}}$$



of inputs coming to the node



$$\text{np.random.randn}(2, 2) * \sqrt{\frac{2}{2}}$$

$$\sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$$

$$\text{He Normal} \quad \text{relu} \quad \left[\frac{1}{\sqrt{n}}, \right)$$

$$\boxed{\sqrt{\frac{2}{\text{fan-in}}}}$$

Uniform Distribution

Xavier Uniform

$[-\text{limit}, \text{limit}]$



$$\text{limit} = \sqrt{\frac{6}{(\text{fan-in} + \text{fan-out})}}$$

He Uniform

$[-\text{limit}, \text{limit}]$

$$\text{limit} = \sqrt{\frac{6}{\text{fan-in}}}$$

Keras

Weight Initialization Techniques

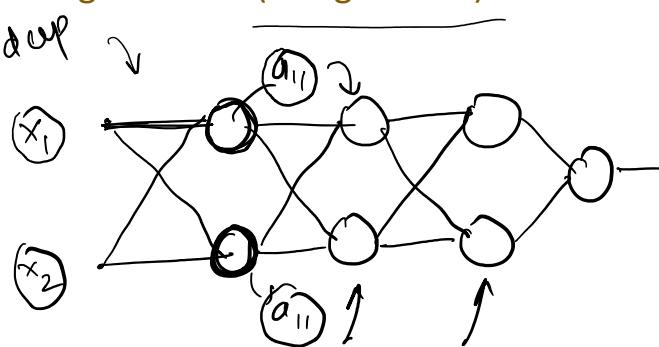
22 June 2022 12:49

What is Batch Norm?

27 June 2022 11:00

Batch-Normalization (BN) is an algorithmic method which makes the training of Deep Neural Networks (DNN) faster and more stable.

- It consists of normalizing activation vectors from hidden layers using the mean and variance of the current batch. This normalization step is applied right before (or right after) the nonlinear function.



normalized

$$\begin{cases} m=0 \\ sd=1 \end{cases}$$

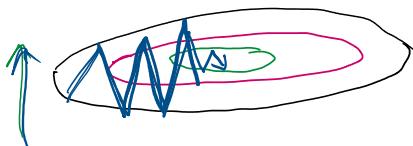
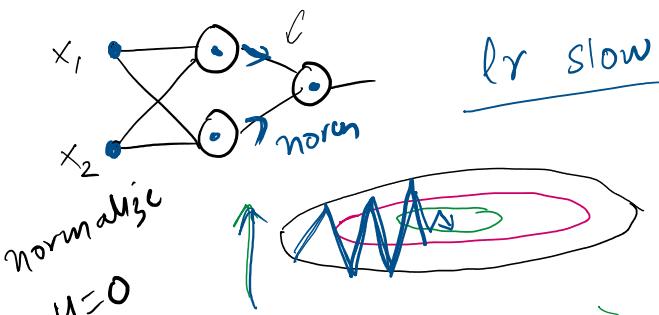
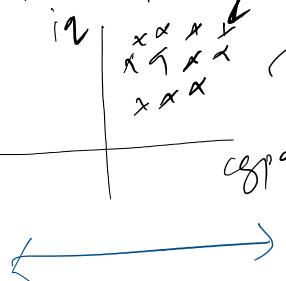
normalize

$$\begin{cases} m=0 \\ sd=1 \end{cases}$$

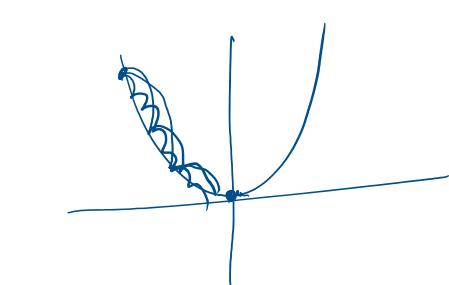
Why use Batch Norm?

30 June 2022 16:08

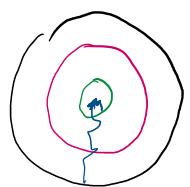
cgpa	iq	placed
7	70	1
8	80	0
9	90	1
6	60	0



lr slow



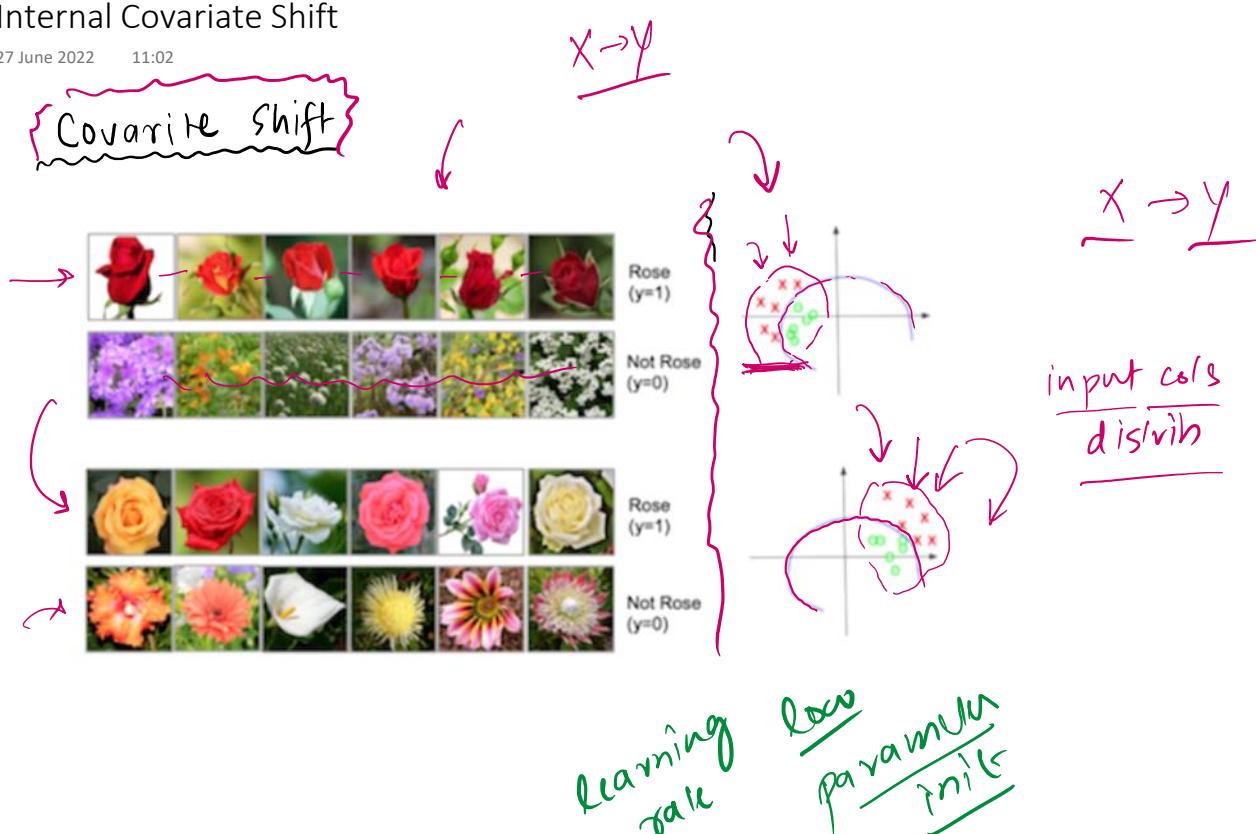
{ training }
slow



training
faster
stable

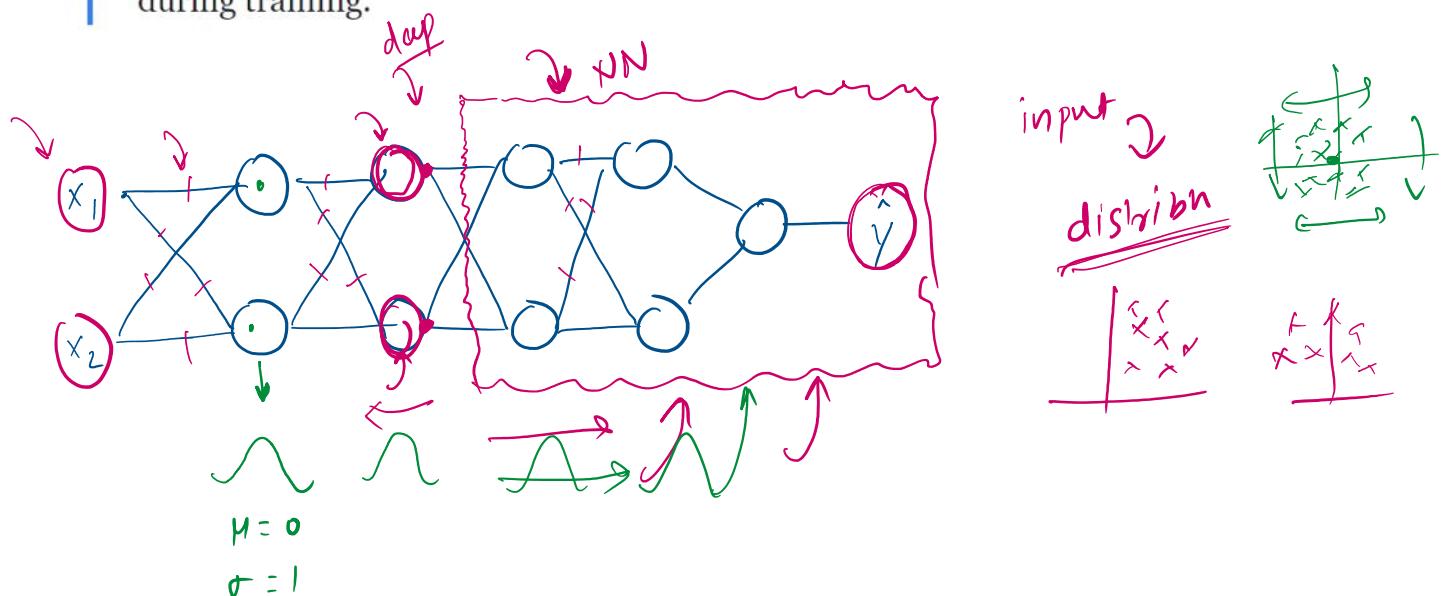
Internal Covariate Shift

27 June 2022 11:02



The authors' precise definition is:

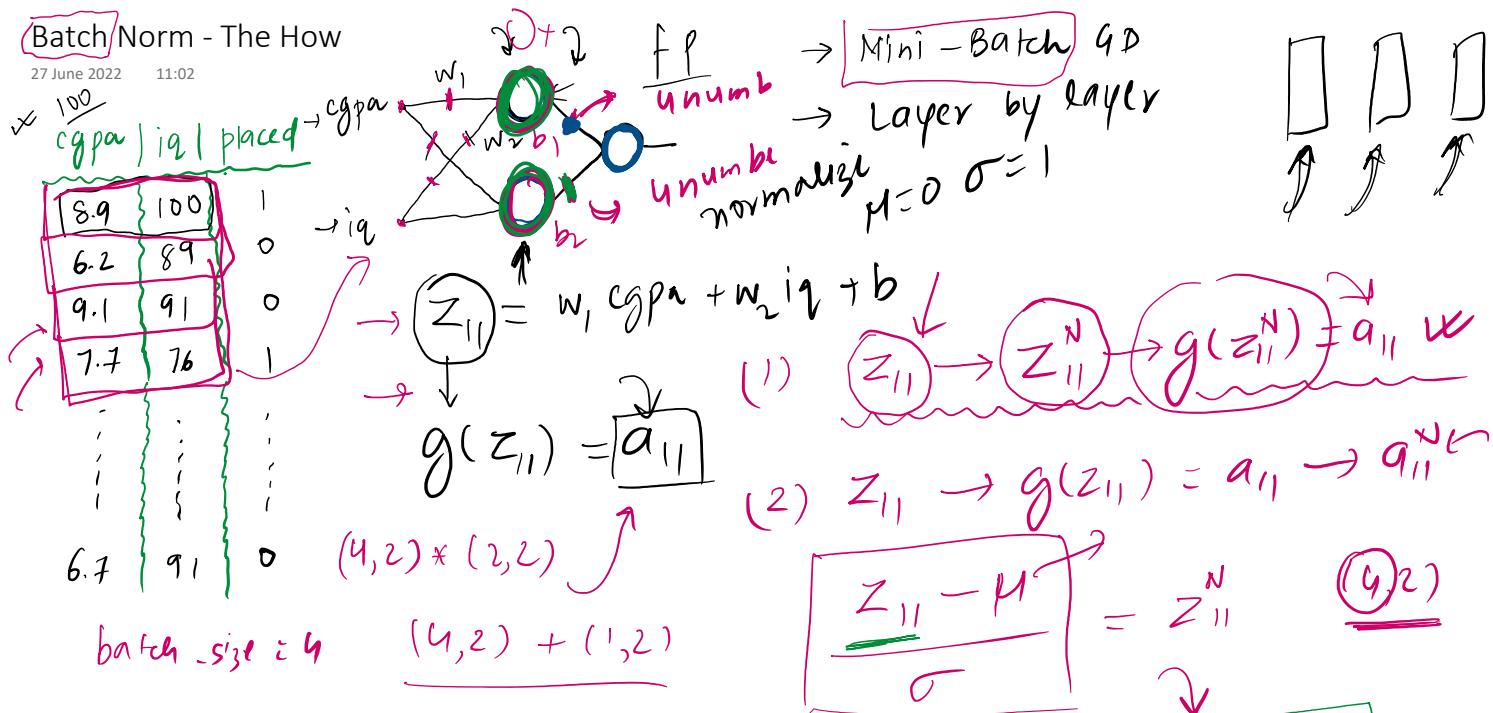
We define Internal Covariate Shift as the change in the distribution of network activations due to the change in network parameters during training.





Batch Norm - The How

27 June 2022 11:02



$$\mu_B = \frac{1}{m} \sum_{i=1}^m z_i$$

$$m = 4$$

$$z_{11}^i = z_{11}^i - \mu_B$$

$$\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_{ii} - \bar{B}_B)^2}$$

$$\boxed{0-1} \quad M=0 \quad \boxed{0-11} \quad \text{learnable parameters} \quad \underline{w, b}$$

$$\boxed{\gamma=1 \quad \beta=0}$$

$$z_{11} \rightarrow z_{11}^N \rightarrow z_{11}^{BN} \rightarrow g(z_{11}^{BN}) = a_{11}$$

~~$\mu = 0$~~ $\sigma = 1$

$$\gamma = \tau + \epsilon \quad \beta = M$$

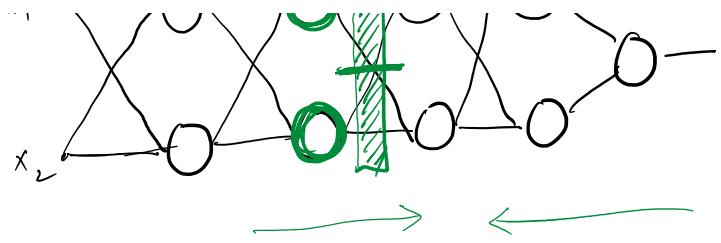
flexibility

$$Z_{11} \rightarrow Z_{11}^N \rightarrow Z_{11}^{BN}$$

Batch Norm → Laym

$$\underline{\gamma} = \underline{\gamma} - \eta \frac{\partial \Phi}{\partial \underline{\gamma}} \quad \text{Keras}$$

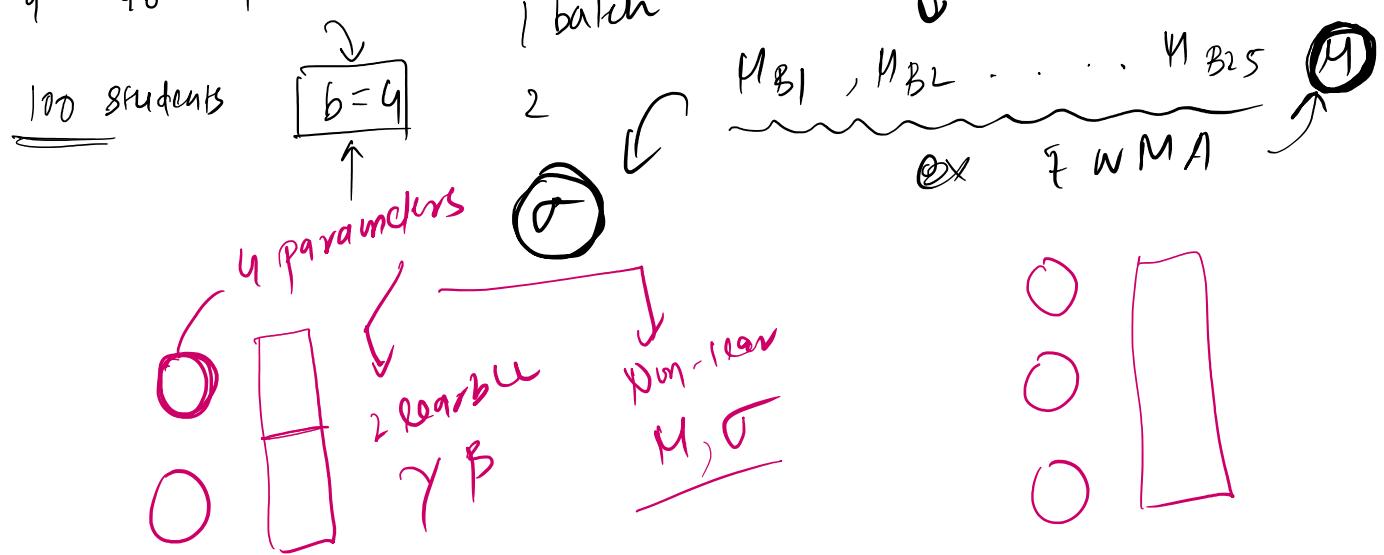
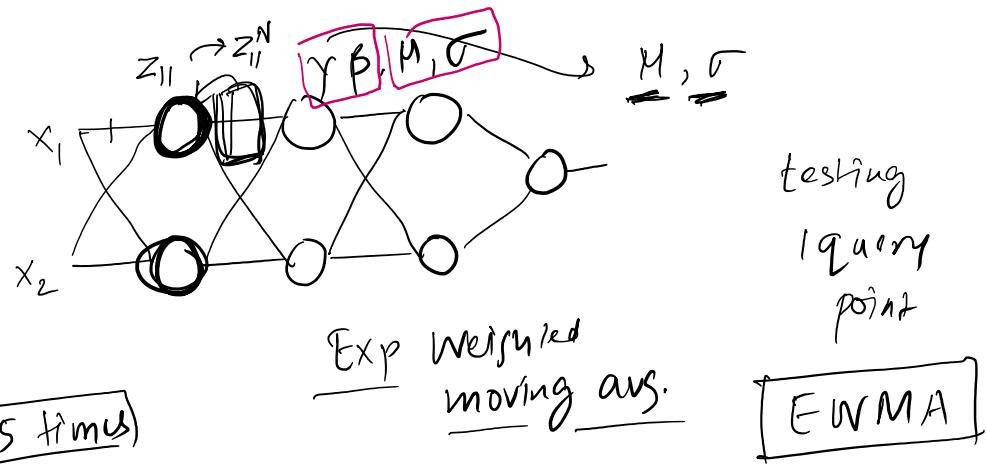
A hand-drawn diagram of a neural network layer. It consists of several nodes connected by lines. One node in the middle row is highlighted with a green circle and has a green cross through it, indicating it is not active. Above the network, the word "flexible" is written in pink with a wavy underline. To the right, the text "Batch N" is written in green. Below the network, the input x_1 is shown at the start of one of the lines. The connections between nodes are labeled with green letters: "u" above the top-left connection, "learnable" above the top-middle connection, " γ_B " above the top-right connection, and " γ_B " above the bottom-right connection.



Batch Norm during test

27 June 2022 11:03

cgpa	ia	placed
8	80	1
7	70	0
6	60	1
:	:	:
9	90	1



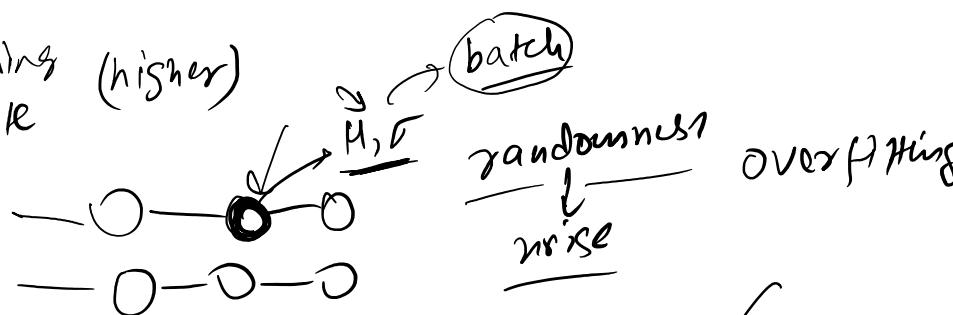
$$3 \times 4 = 12$$

6

6

Advantages

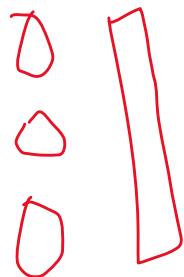
27 June 2022 11:00

- 1) stable → hyperpara → wider range of robustness
- 2) faster → learning rate (higher)
batch
- 3) Regularizer →
↓
dropout

- 4) weight init impact reduce


Keras Implementation

27 June 2022 11:03

```
model = Sequential()  
  
model.add(Dense(3,activation='relu',input_dim=2))  
model.add(BatchNormalization()) ←  
model.add(Dense(2,activation='relu')) ←  
model.add(BatchNormalization()) ←  
model.add(Dense(1,activation='sigmoid'))
```

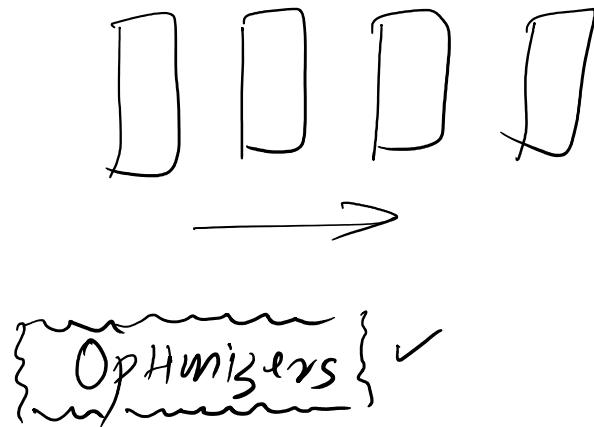


Introduction

05 July 2022 09:57

Performance of NN

↓ how the speed the
training

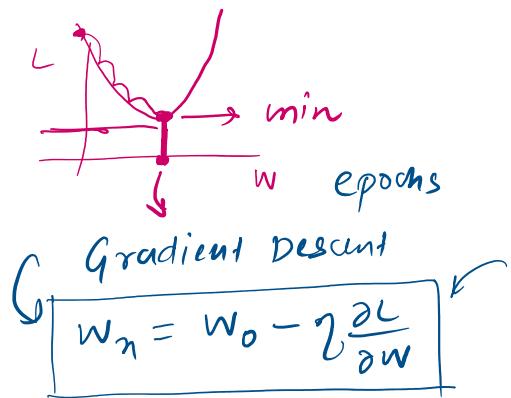
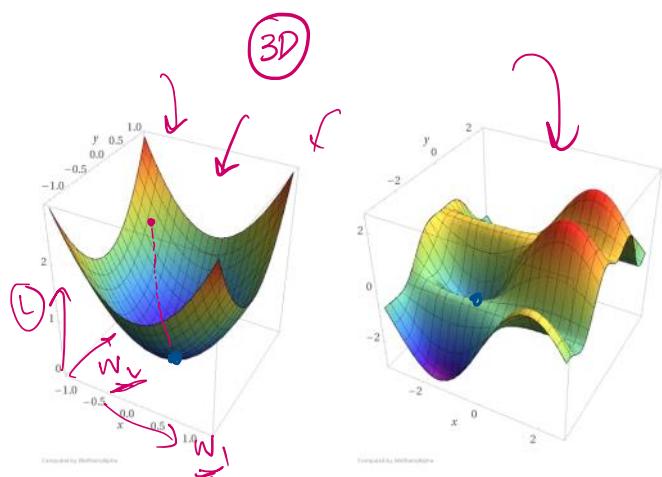
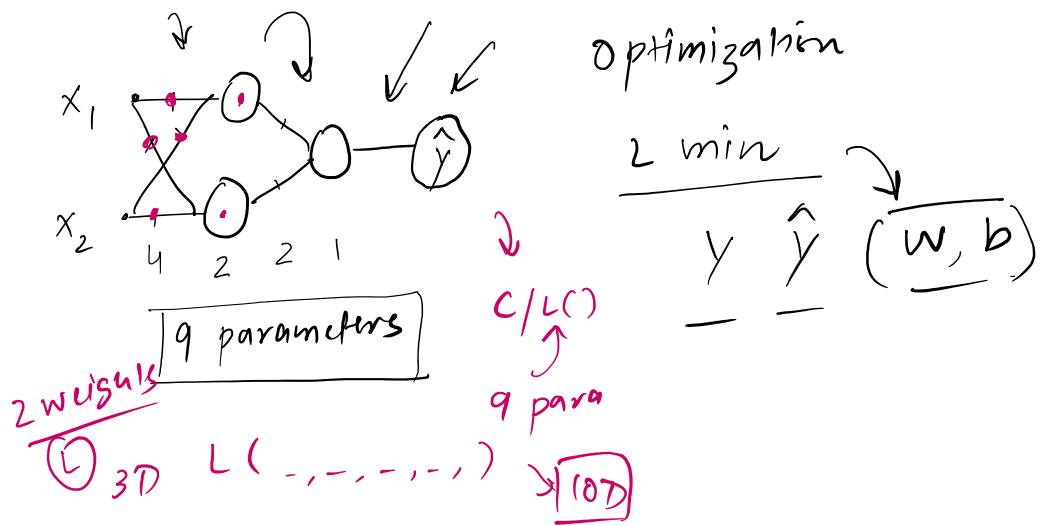


- Weigh init
- Batch Norm
- Activation function

Role of Optimizer

05 July 2022 10:01

cgpa		ia	placed
8	80	1	1
9	90	1	0
7	70	1	0



Types of Optimizers

05 July 2022 10:02

- Batch GD
- Stochastic GD
- Mini batch GD

epochs = 10

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

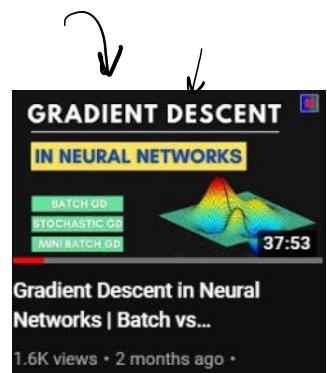
500 rows → pred

↓
loss → weight update

10×500 → 5000

batch size = 100

5 batch → 10×5 times

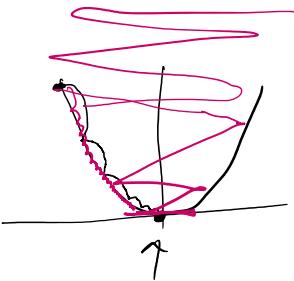


Challenges

05 July 2022 10:02

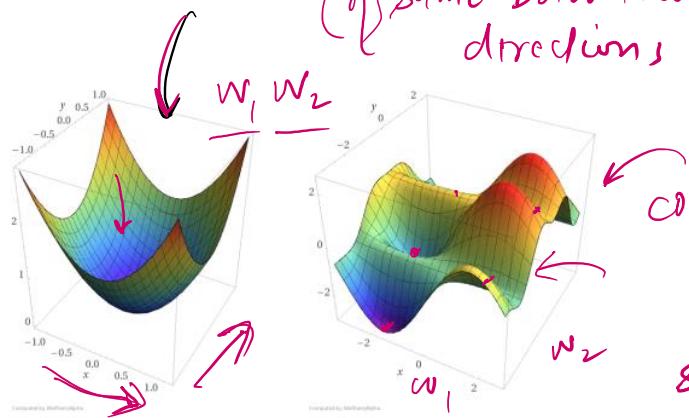
1) learning rate

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



2) learning rate scheduling → pre define

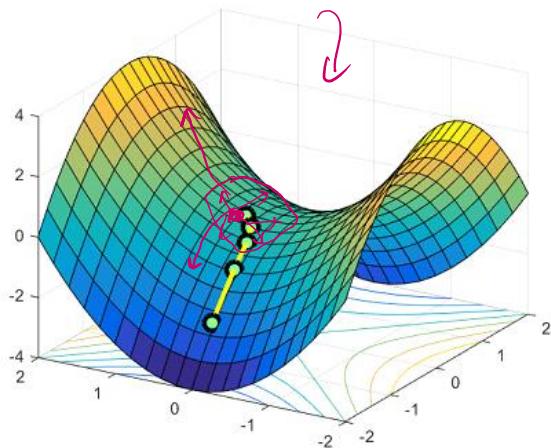
3)



(9) weight's and bias's 9 sep learn
 $\overline{10-D} \rightarrow 10 \text{ direction}$
 complex minima
 many sub optimum

4) local minima

5) saddle point



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

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

What next?

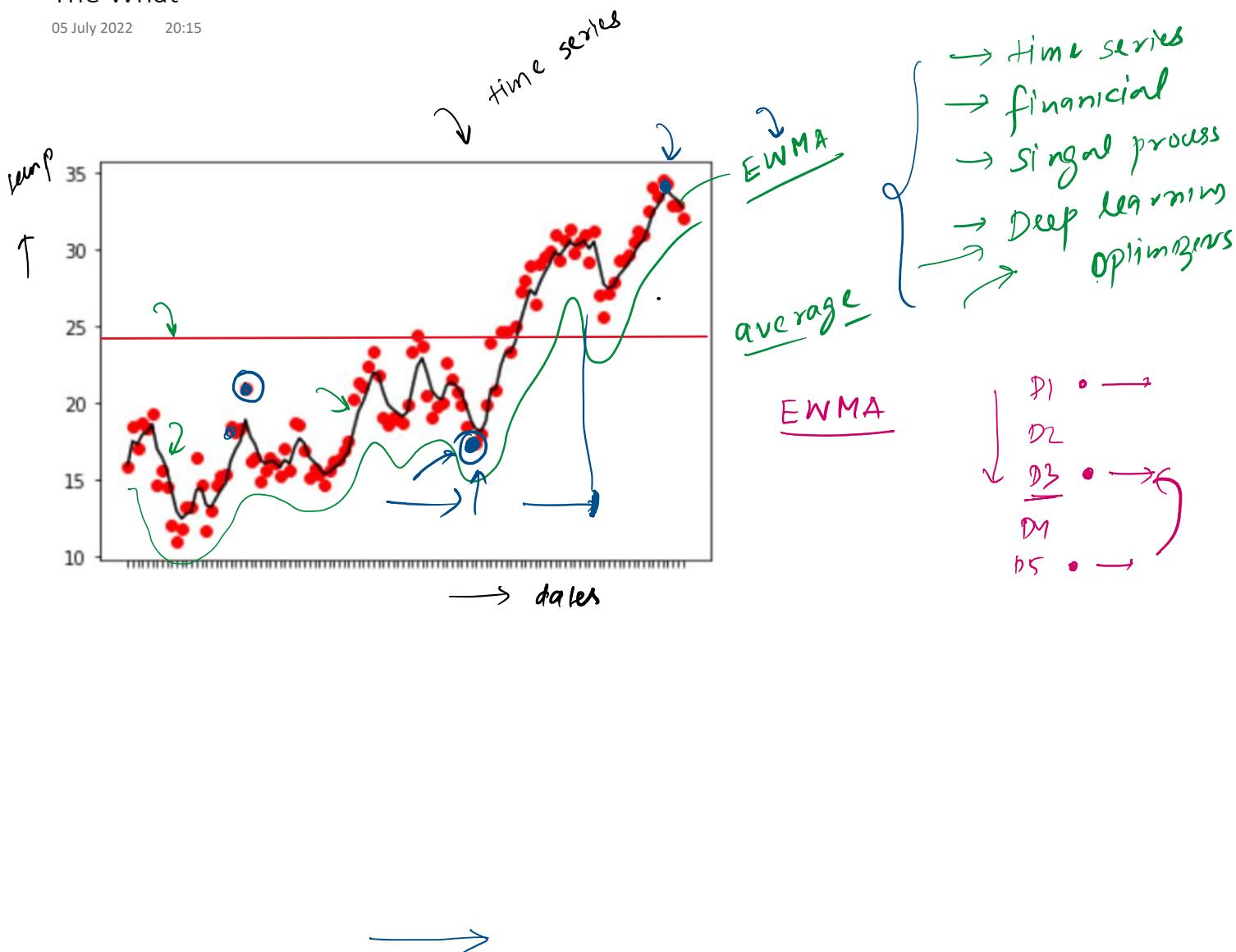
05 July 2022 10:02

- 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

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

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

$$V_0 = 0$$

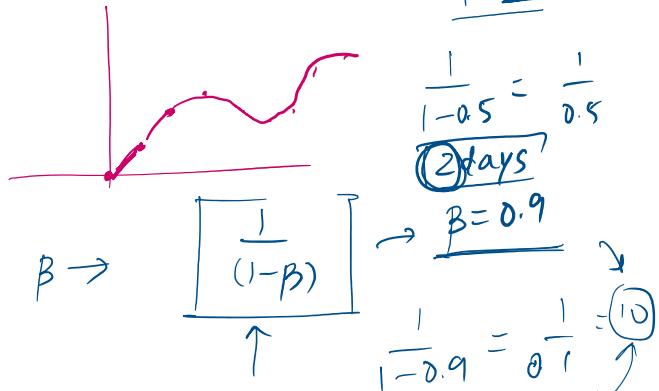
$$V_0 = \theta_0$$

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

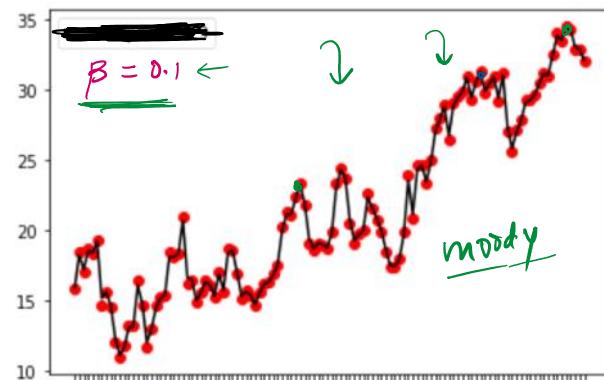
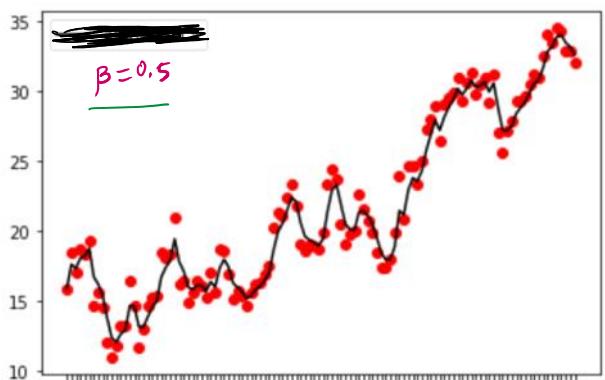
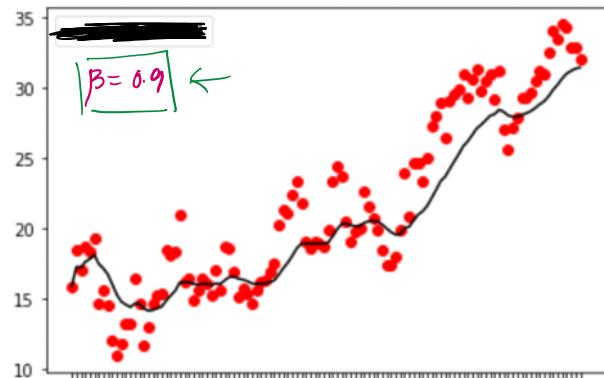
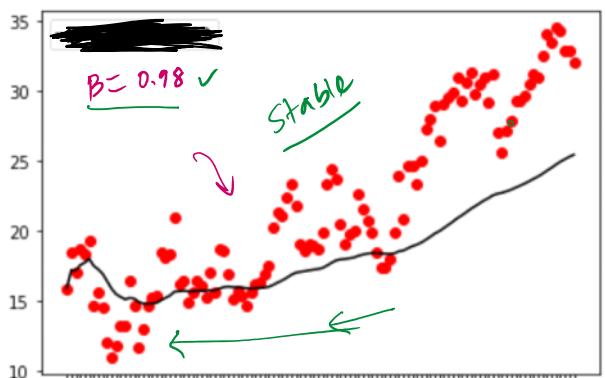
$$V_2 = 0.9 \times 1.3 + 0.1 \times 17 =$$

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

$$\beta = 0.5$$



Effect of β



Mathematical Intuition

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

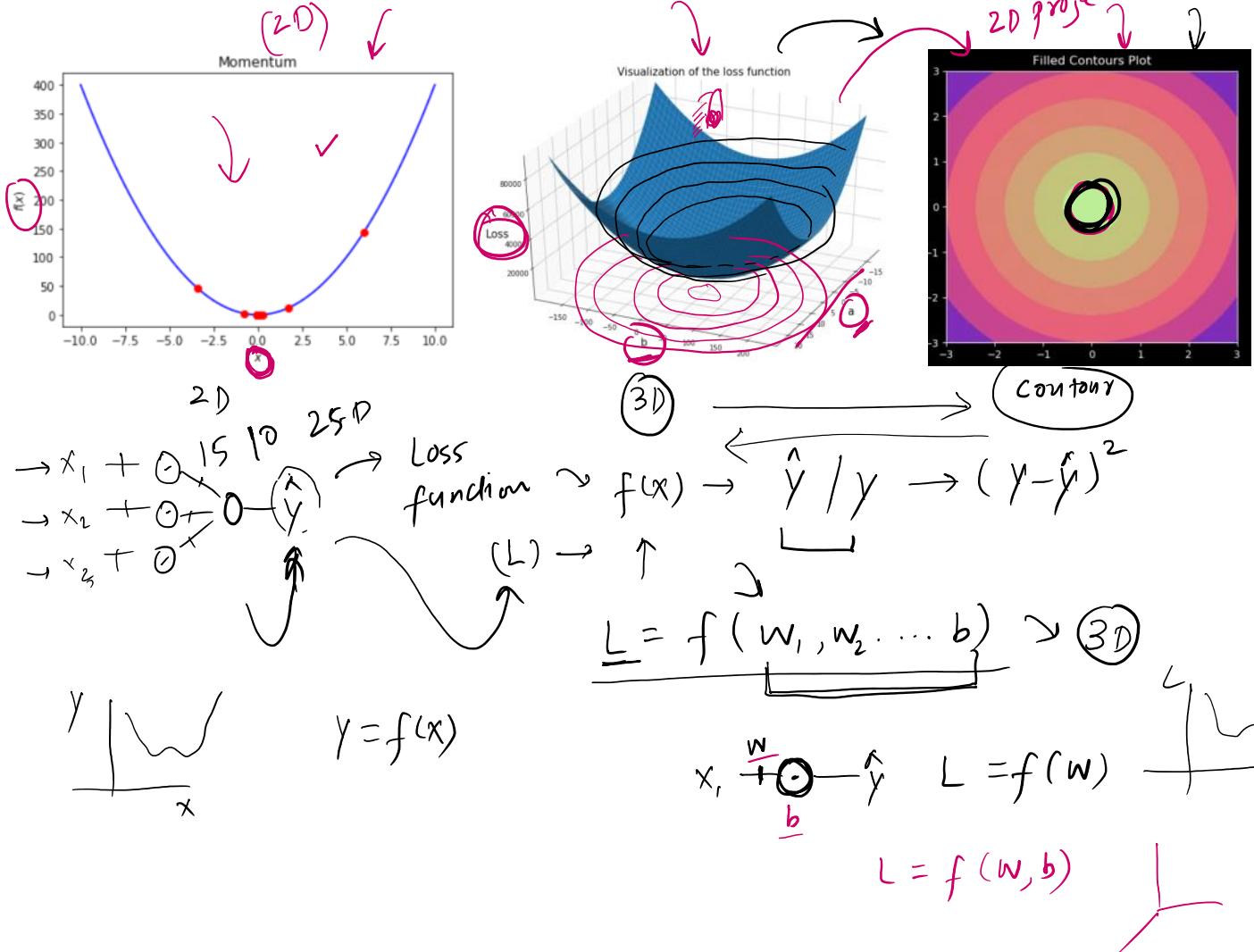
$$\left. \begin{aligned} V_4 &= \beta V_3 + (1-\beta) \theta_4 \\ V_4 &= \beta^3 (1-\beta) \theta_1 + \beta^2 (1-\beta) \theta_2 + \beta (1-\beta) \theta_3 + (1-\beta) \theta_4 \end{aligned} \right\}$$

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

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

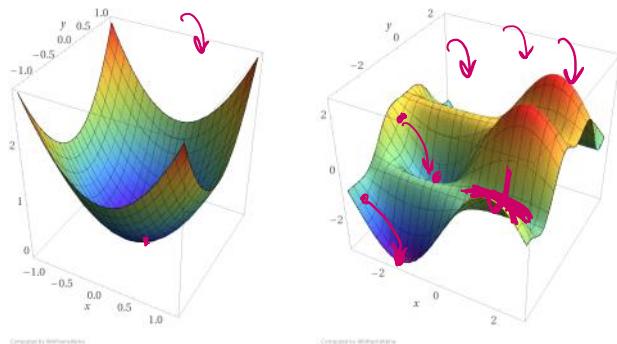
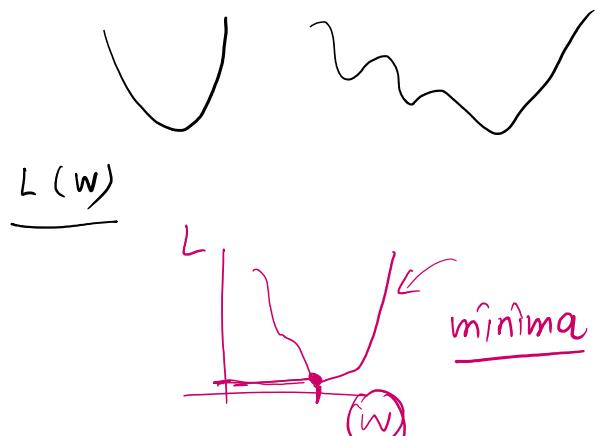
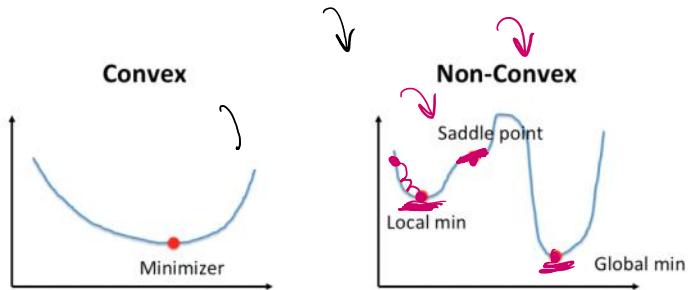
Understanding Graphs

20 July 2022 12:03

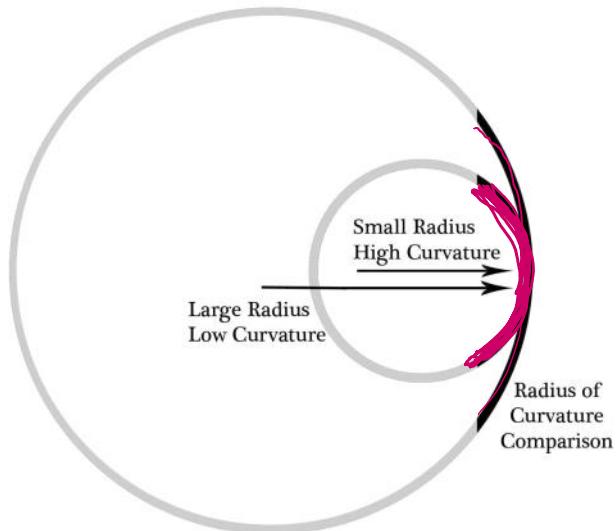


Convex Vs Non-Convex Optimization

20 July 2022 13:06

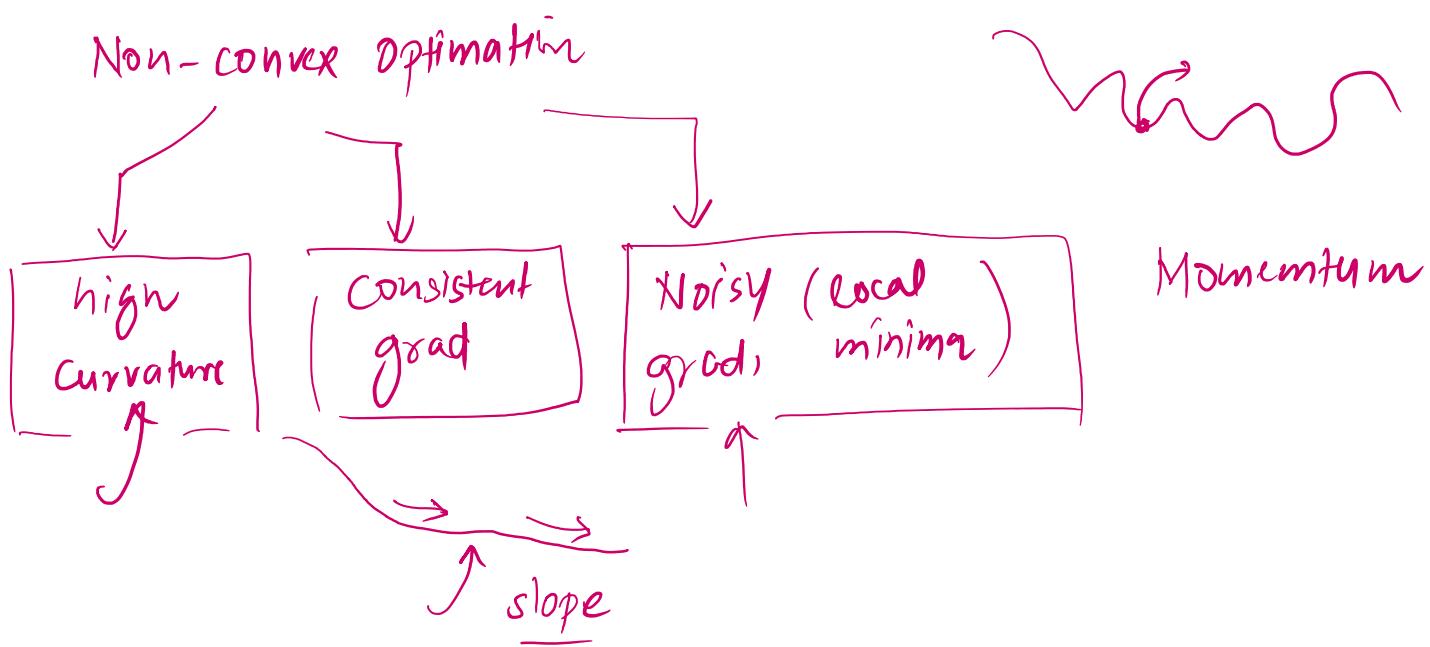


- 1) local minimum
 - 2) saddle point
 - 3) High curvature
-



Momentum Optimization - The Why?

20 July 2022 14:28

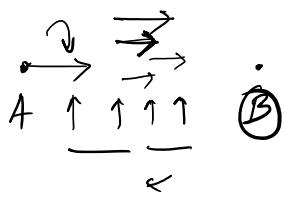


Momentum Optimization - The What?

20 July 2022 14:33

Speed

Common



prev
grad

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

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



Momentum Optimization - Mathematics(The How?)

20 July 2022 14:56

$$\underline{w_{t+1}} = \overrightarrow{\underline{w_t}} - \boxed{\eta \nabla w_t}$$

acceleration

momentum

history of past velo- use

0.9

$v \rightarrow \text{velocity}$

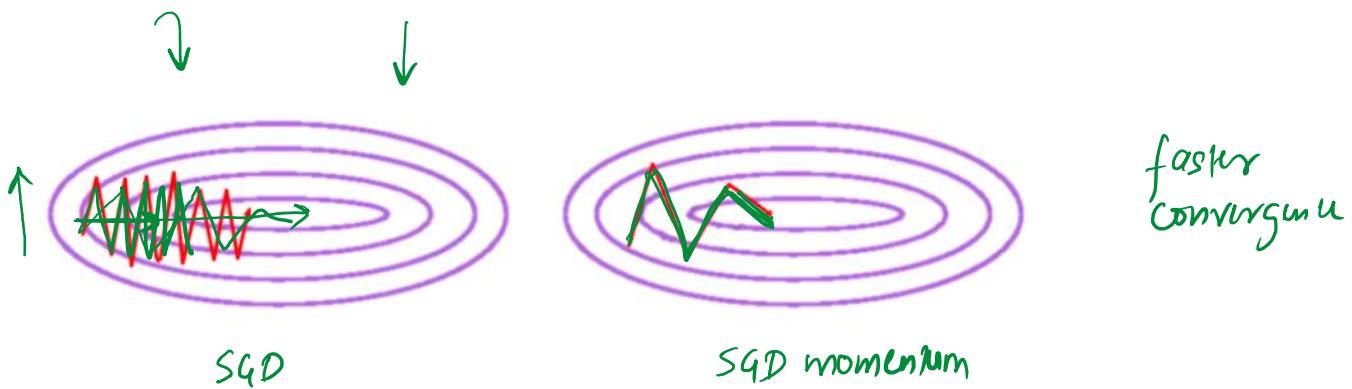
$$w_{t+1} = w_t - v_t \rightarrow$$

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

$0 < \beta < 1$

prev velocity

v_{t-2}
 v_{t-3}



Effect of beta

20 July 2022 15:09

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

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

$\boxed{\beta=0}$

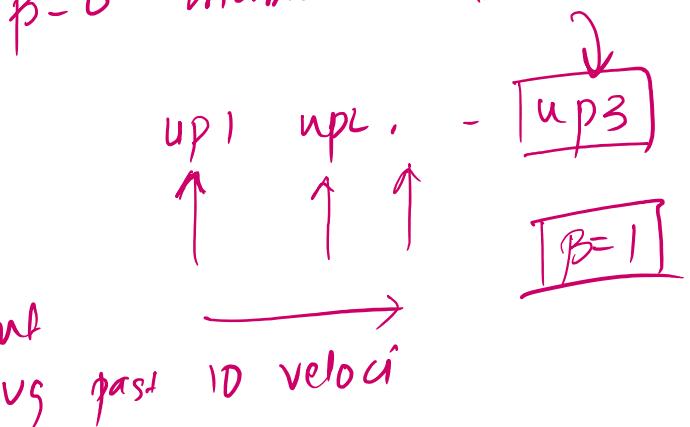
SGD

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

$\beta=0$ momh \rightarrow SGD

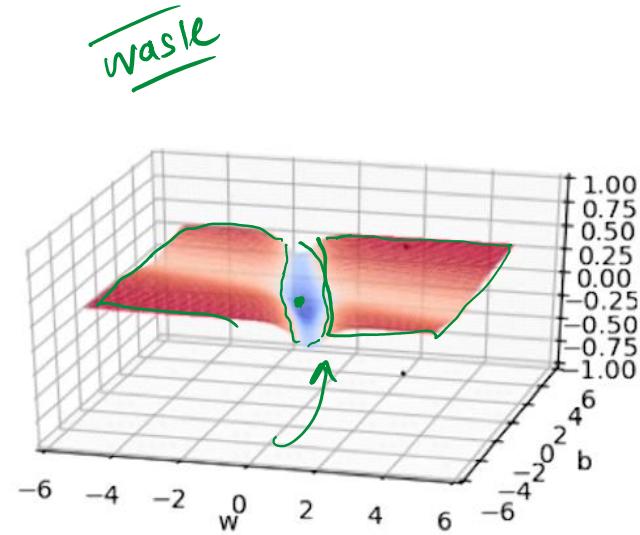
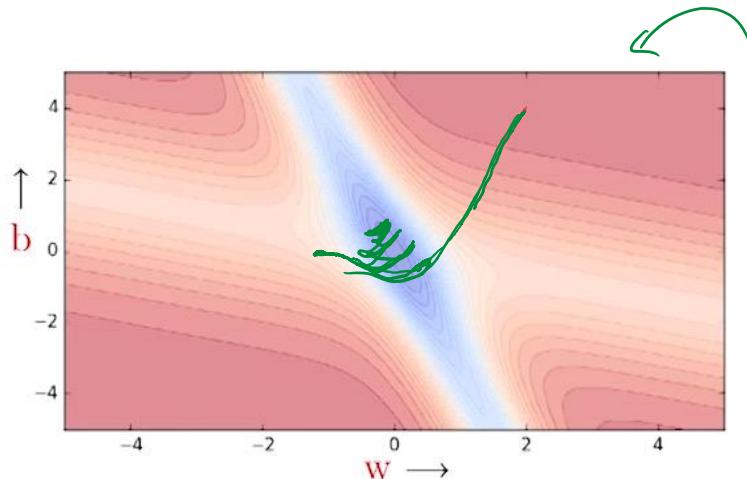
$$\frac{1}{1-\beta} = \frac{1-10}{0.1} \text{ decaying factor}$$

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



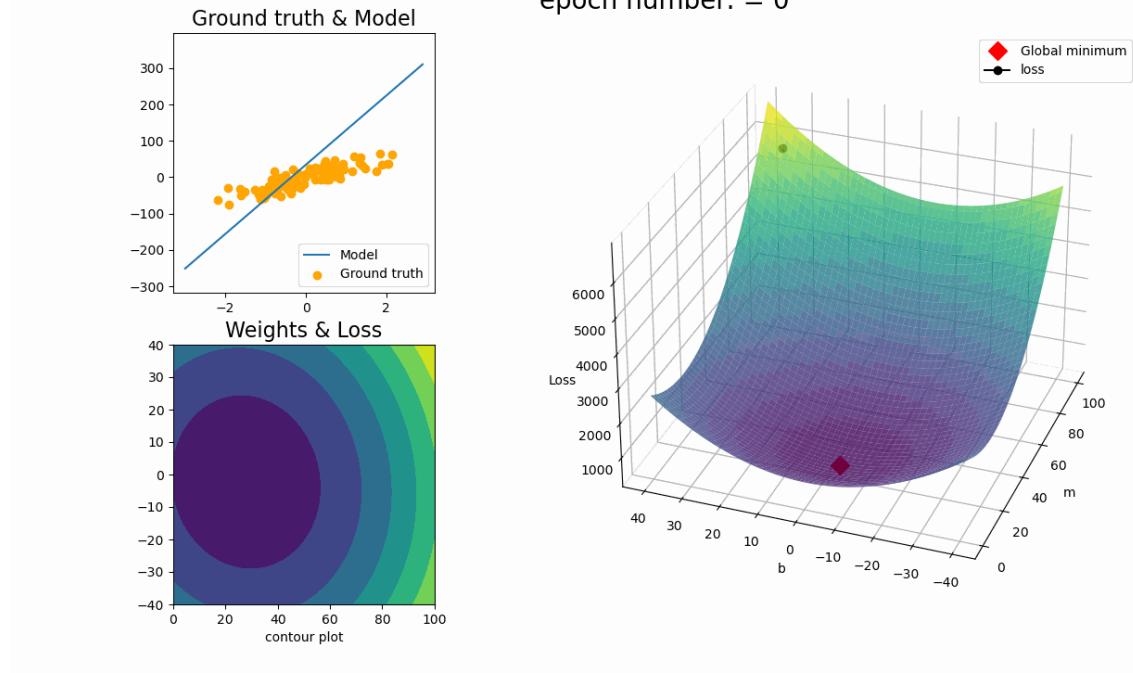
Problem with momentum optimization

20 July 2022 15:16



decay 0.9

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



Demo

24 July 2022 13:17

Mathematical Intuition

24 July 2022 13:17

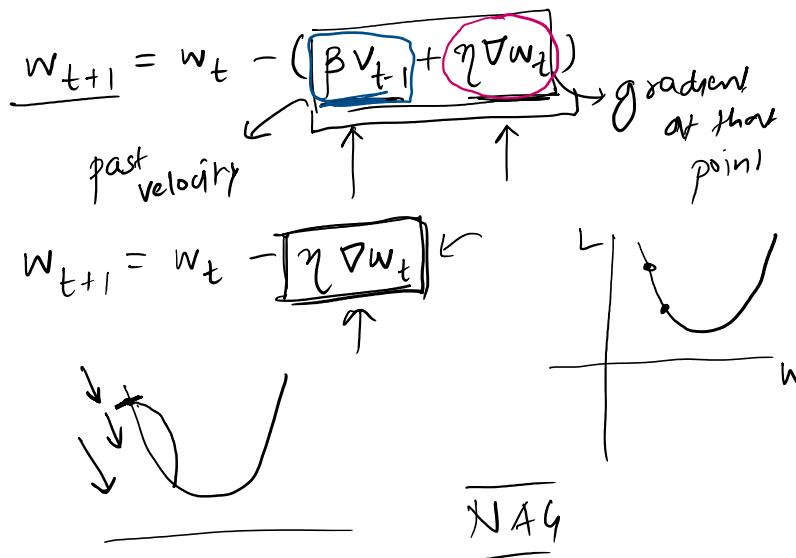
Momentum

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

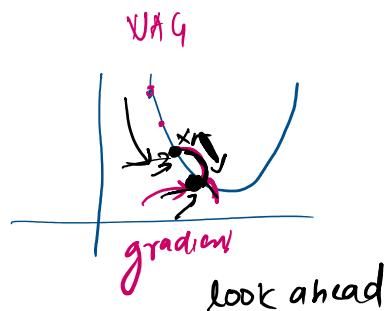
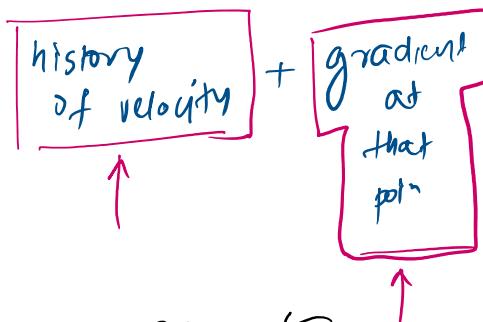
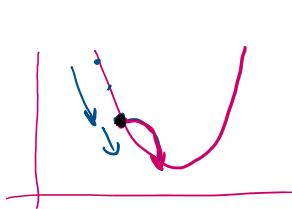
where

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

$\beta \rightarrow$ decay factor



Nesterov Accelerated Gradient *better*

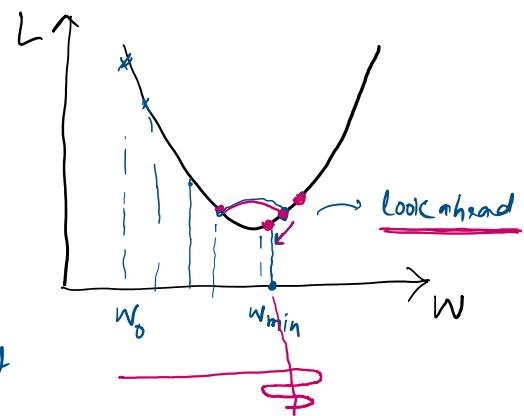
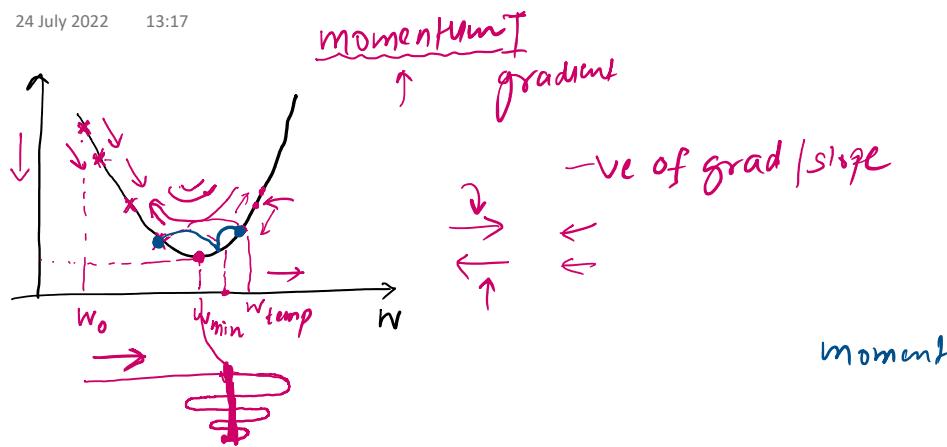


$$\left\{ \begin{array}{l} w_{\text{la}} = w_t - \beta v_{t-1} \\ v_t = \beta v_{t-1} + \eta \nabla w_{\text{la}} \\ w_{t+1} = w_t - v_t \end{array} \right.$$

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

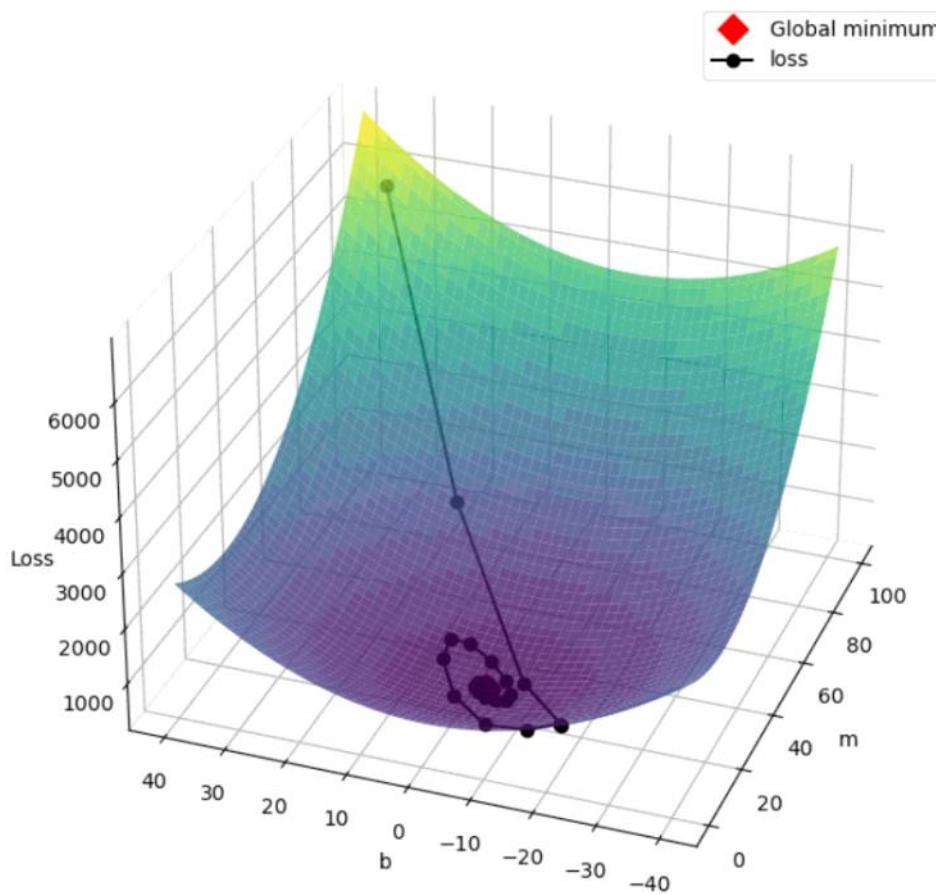
Geometric Intuition

24 July 2022 13:17



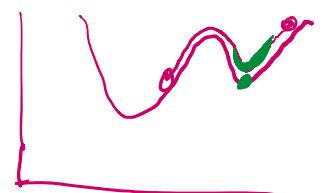
Disadvantage

24 July 2022 13:18



Disadvantage

NAG \rightsquigarrow oscillation
 \downarrow
dampen
local minimum



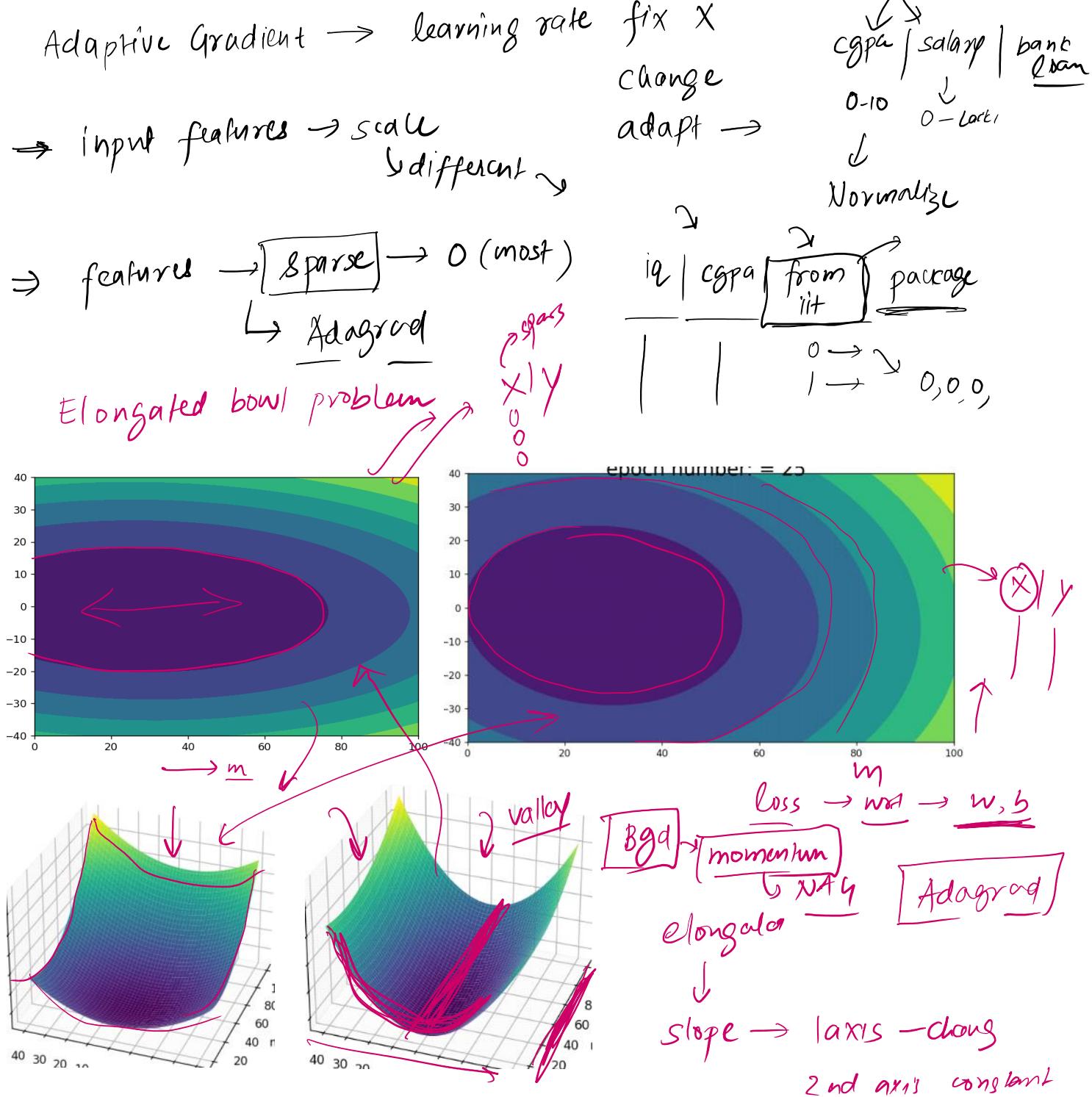
Keras Code

24 July 2022 13:18

2

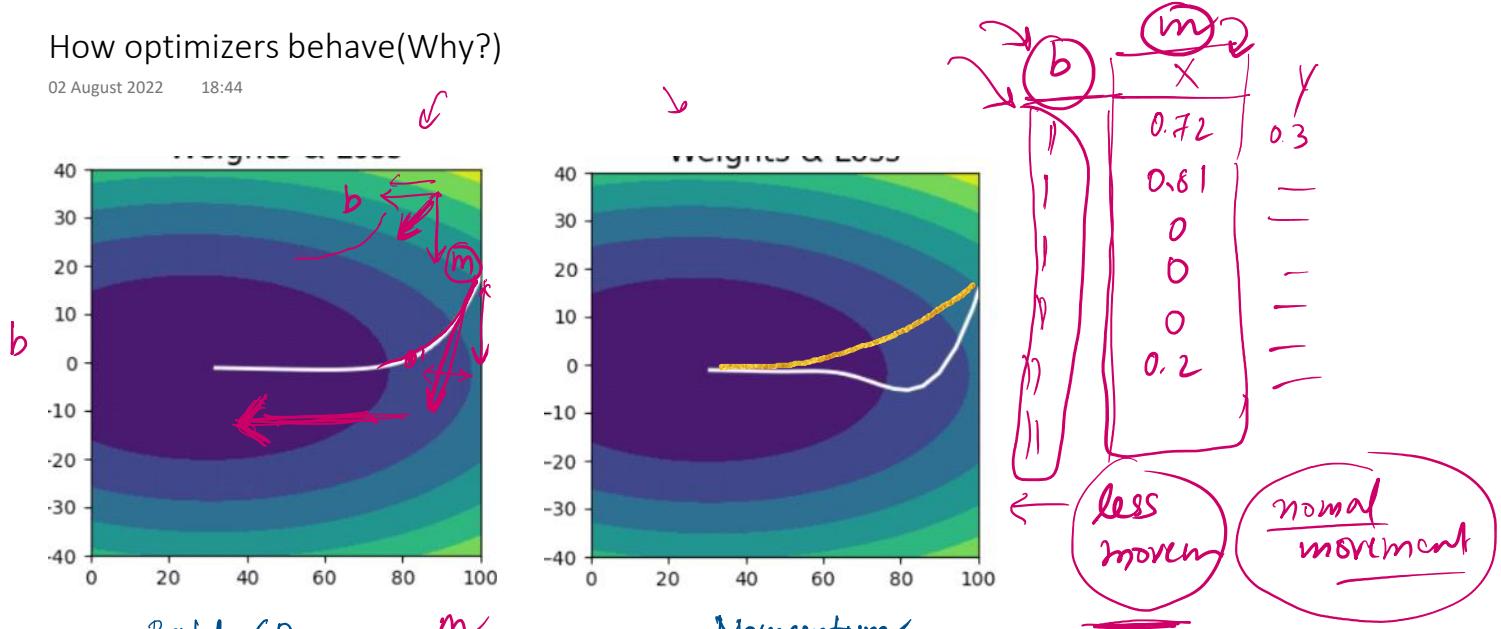
```
tf.keras.optimizers.SGD(  
    learning_rate=0.01, momentum=0.0, nesterov=False, name="SGD", **kwargs  
)
```

SGD } Momentum } NAG
 | momentum = 0.9 | momentum = 0.9
 | nesterov = False | nesterov = True
 | ✓ | ✓



How optimizers behave(Why?)

02 August 2022 18:44



$y \hat{y}$

$x w b \hat{y}$

linear

$\frac{\partial L}{\partial w} \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w}$

$= -2(y - \hat{y})^T X$ add

$\frac{\partial L}{\partial b} = \boxed{-2(y - \hat{y})^T \mathbb{1}}$ $b \leftarrow \text{ignum}$

$BGD \rightarrow \text{sparse}$

for i in epochs:

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

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

100 rows \downarrow small in every epoch

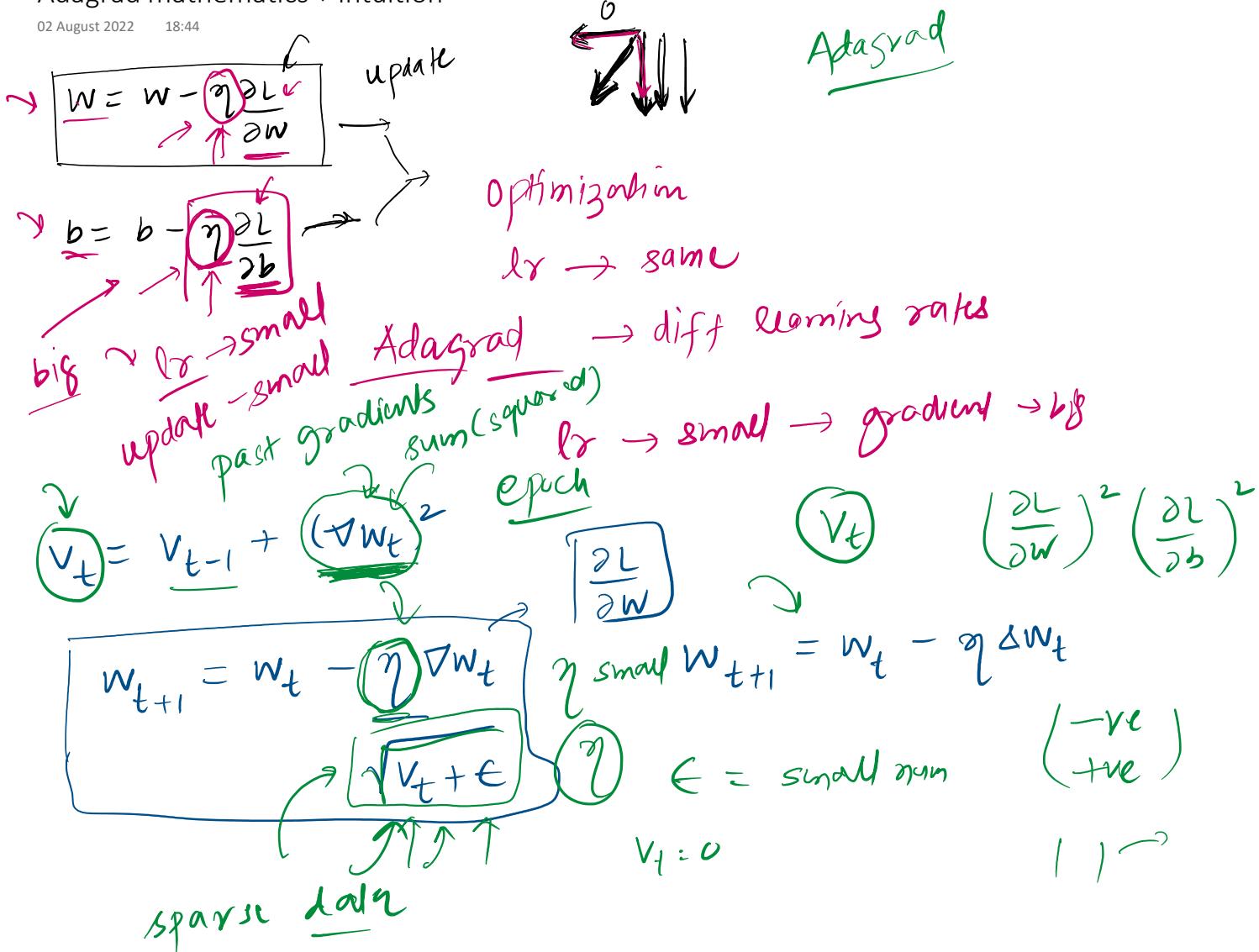
$\times \text{sparse}$

\rightarrow big update in every epoch

sparse col \rightarrow small non sparse
normal/big

Adagrad mathematics + Intuition

02 August 2022 18:44



Adagrad Demo

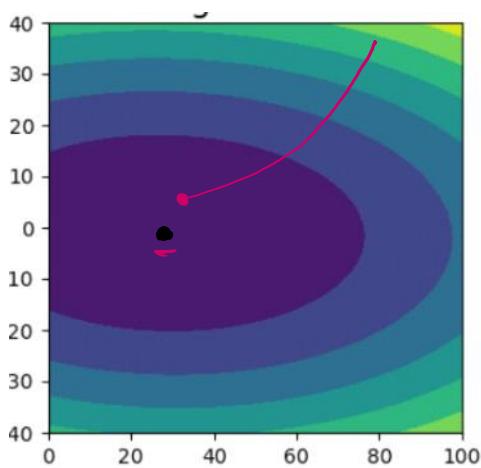
02 August 2022 18:44

Keras Implementation

02 August 2022 18:44

Disadvantage

02 August 2022 18:43



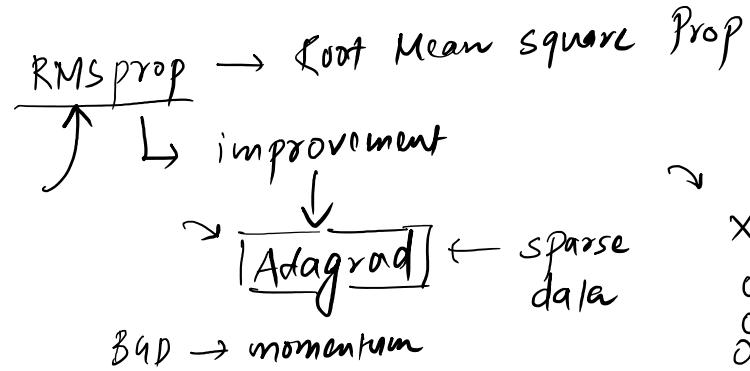
NN \rightarrow use $\times \curvearrowright$ linear regress

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

η $\xrightarrow{\text{small}}$ slower
 $\sqrt{v_t}$ $\xrightarrow{\text{large}} \xrightarrow{\text{past updates}} \xrightarrow{\text{gradient}} \xrightarrow{\text{epochs}}$
 v_t $\xrightarrow{\text{small}} \xrightarrow{\text{update}} \xrightarrow{\text{global minima}} \xrightarrow{\text{RMSprop}} \xrightarrow{\text{Adam}}$

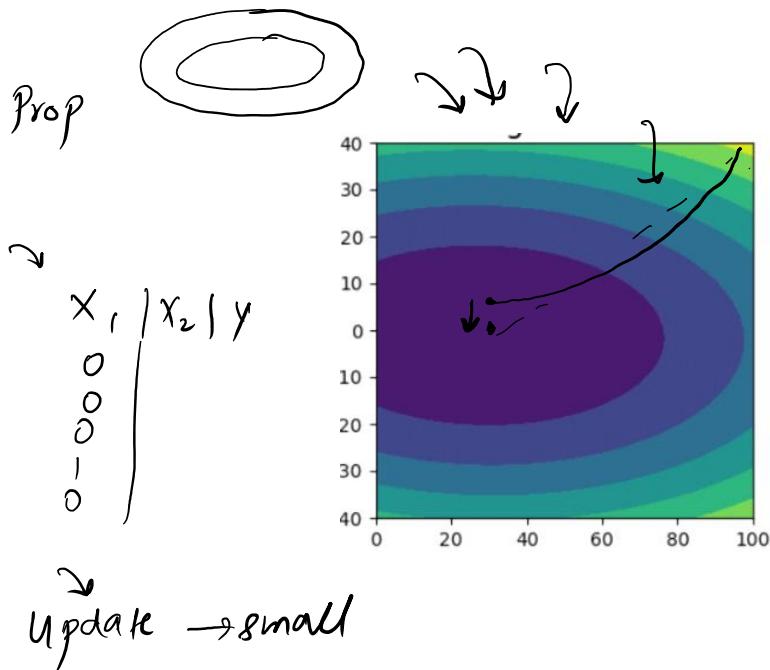
The Why?

03 August 2022 14:03



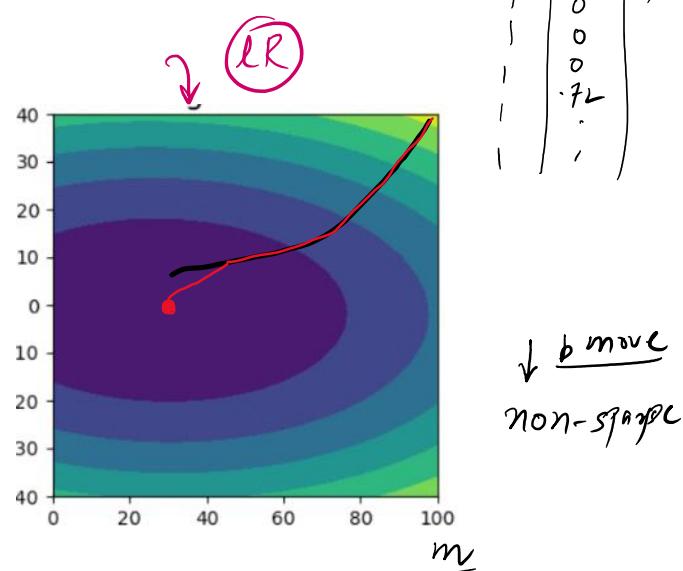
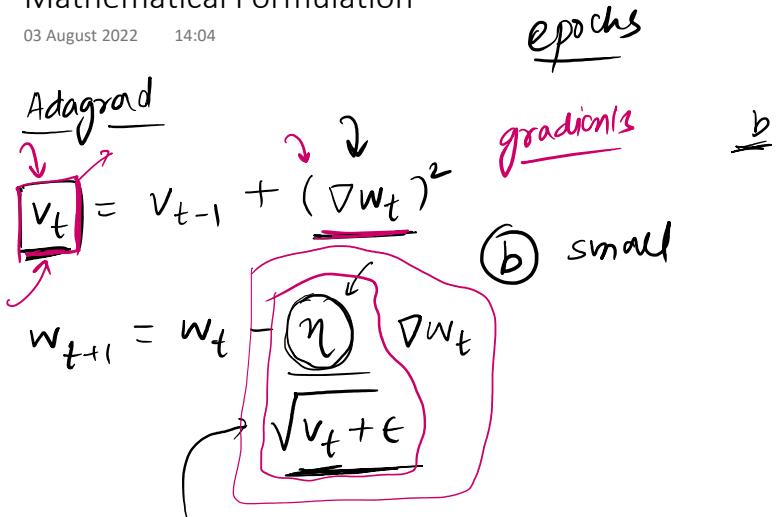
Disadvantage → learning rate

↓



Mathematical Formulation

03 August 2022 14:04



RMSprop

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

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

exp decaying avg
old epoch
 $\beta = 0.95$

$$\rightarrow v_0 = 0$$

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

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

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

larger

\uparrow smaller \uparrow epoch 1 \uparrow epoch 2 \uparrow epoch 3

v_t shoot $\rightarrow \eta$ small

Adagrad \rightarrow NN \rightarrow non-convex optimization

convex optimis \rightarrow linear $\frac{\partial \varphi}{\partial \varphi}$

✓ convex optimis \rightarrow unconstrained

global minima RMSprop

Disadvantage

03 August 2022

14:54

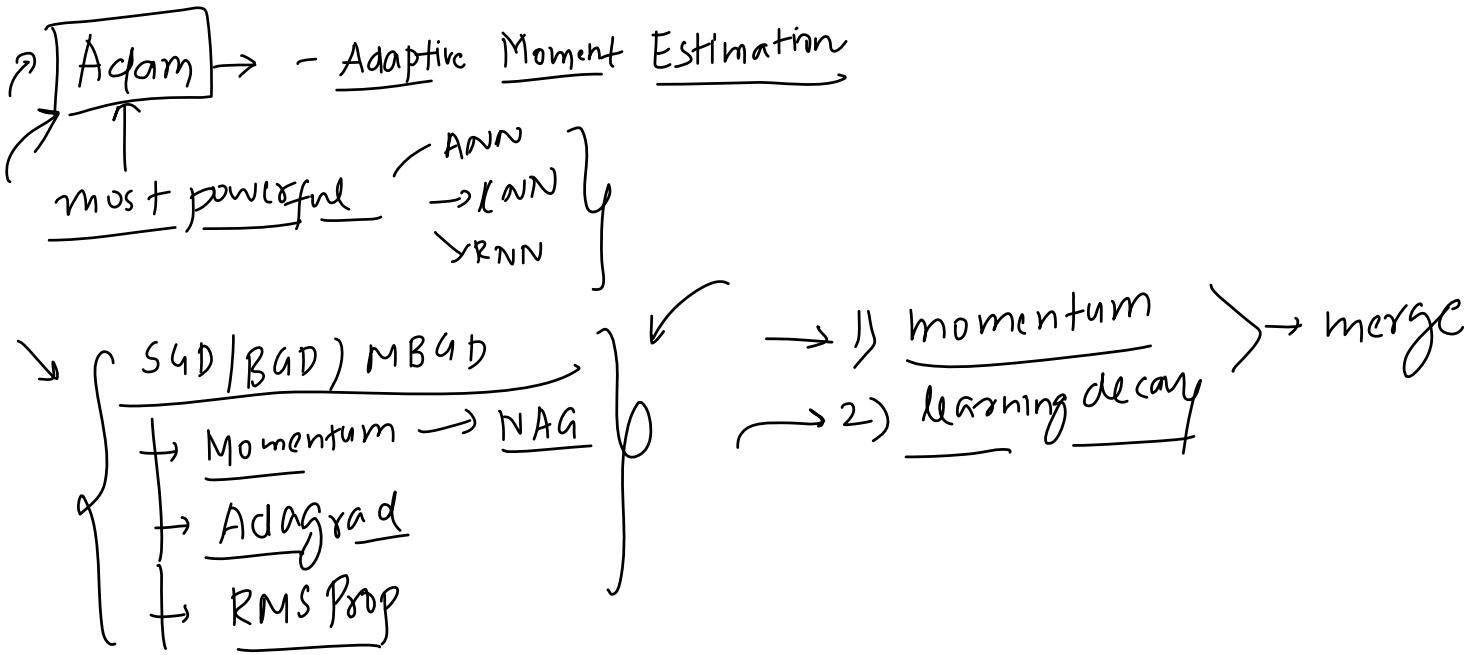
No,

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

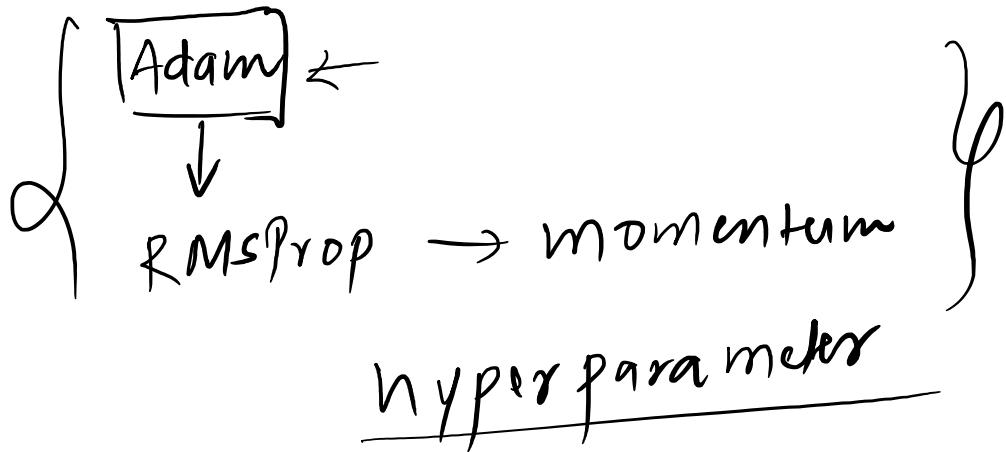
$m_0 = 0 \quad v_0 = 0 \quad \eta = 0.1 \quad \beta_1 = 0.9 \quad \beta_2 = 0.99$

Demo

04 August 2022 10:45

Verdict

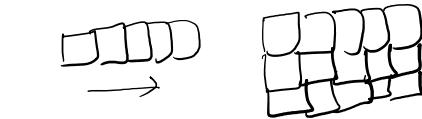
04 August 2022 10:45



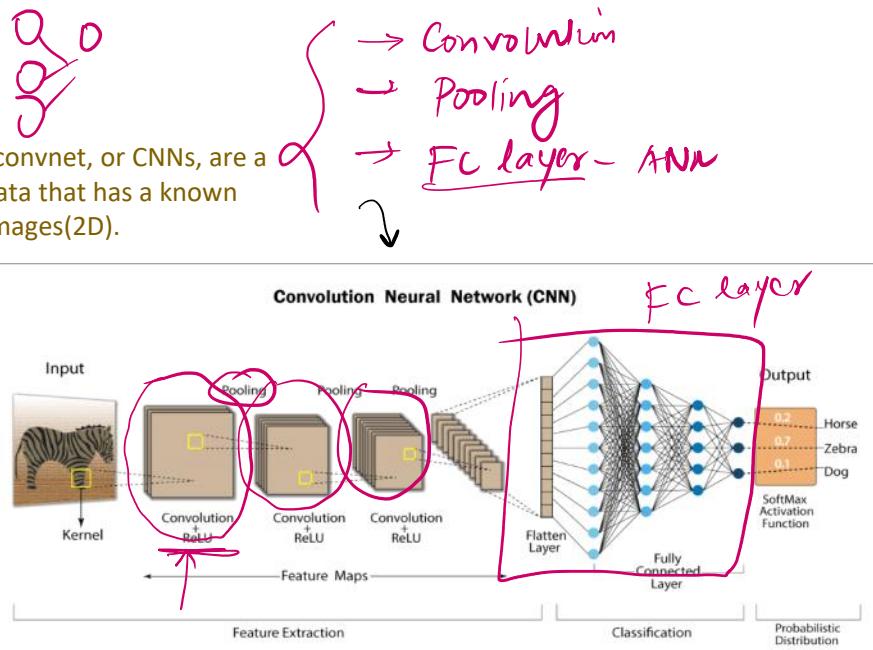
What is a CNN?

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).



→ ANN → Convolution
ANN → matrix mult.



Inspirations

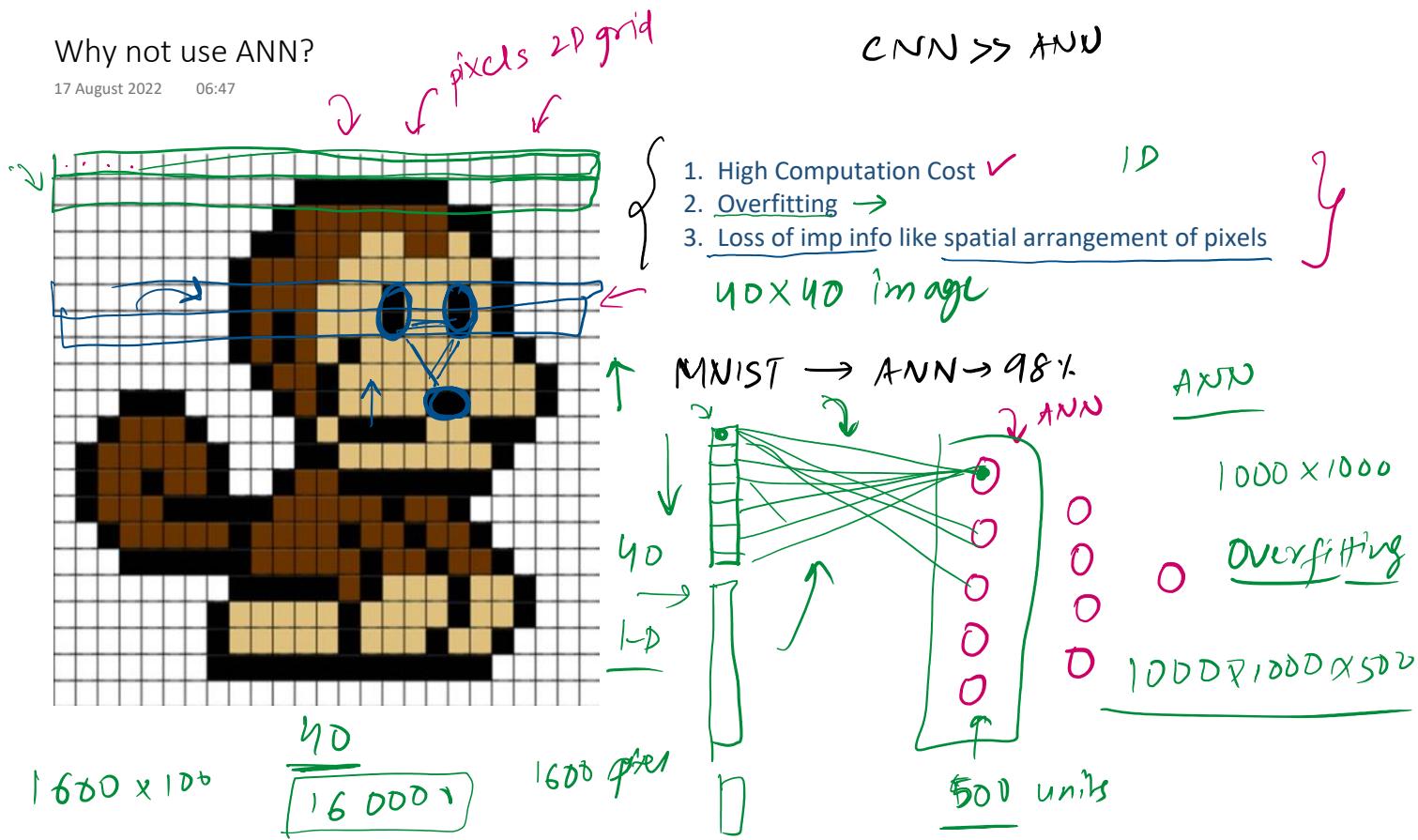
↳ visual cortex

1998 Yann LeCun → AT&T →
↳ Microsoft → OCR hand writing ↳ CNN

↳ facial recog ↳ CNN
self driving ↳ RNN

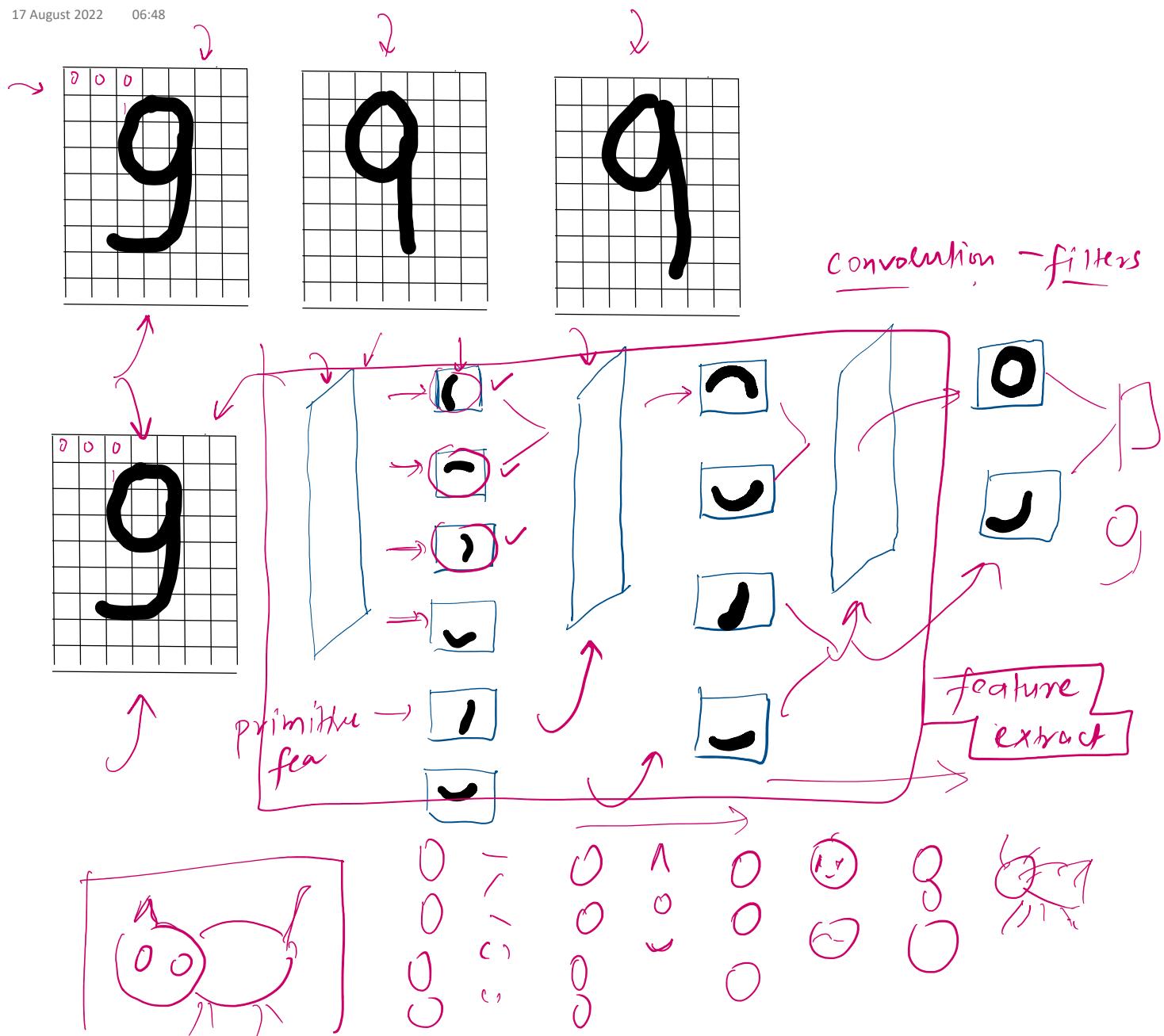
Why not use ANN?

17 August 2022 06:47



CNN Intuition

17 August 2022 06:48

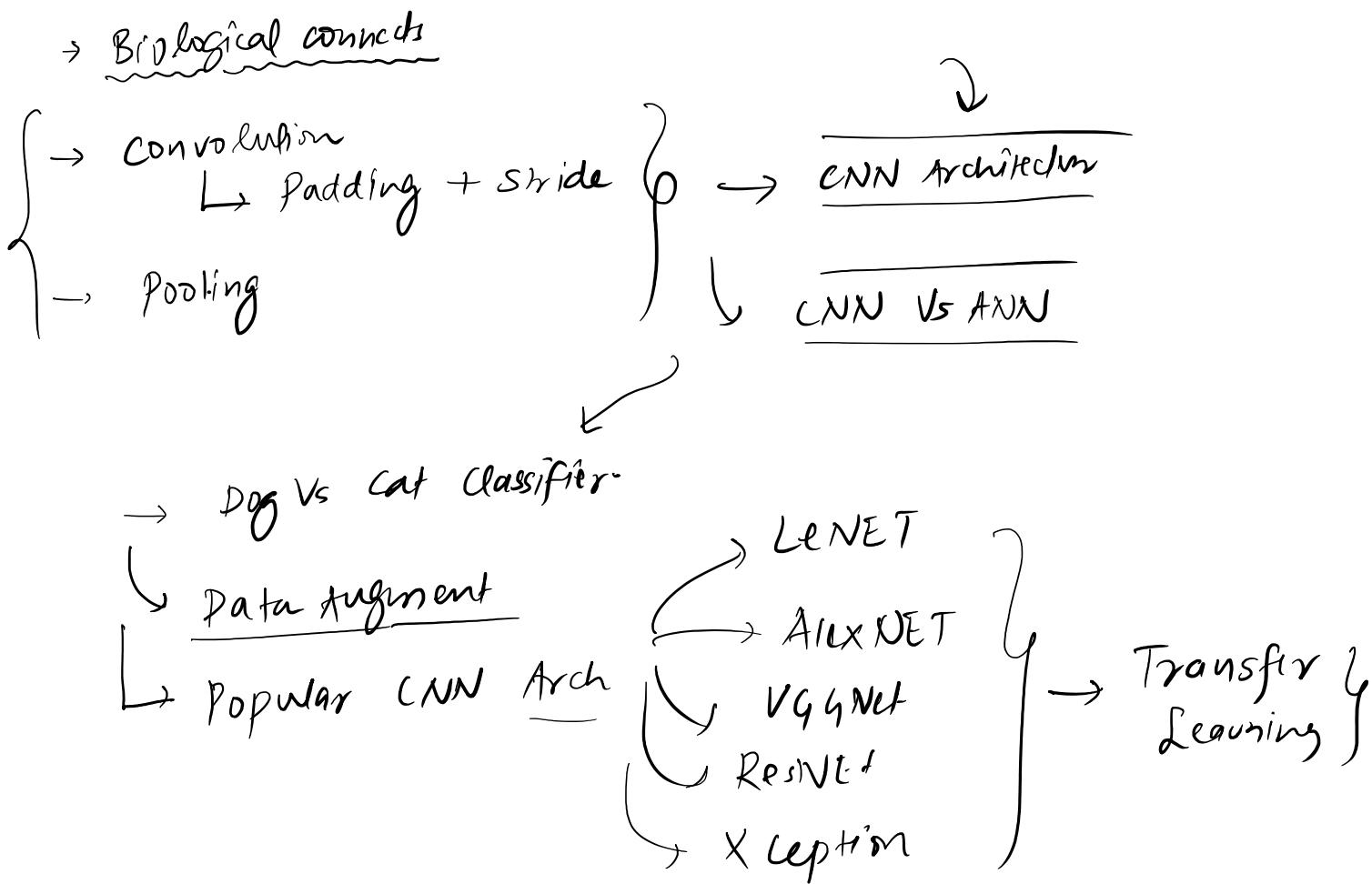


CNN Applications

17 August 2022 06:48

Roadmap

17 August 2022 06:48

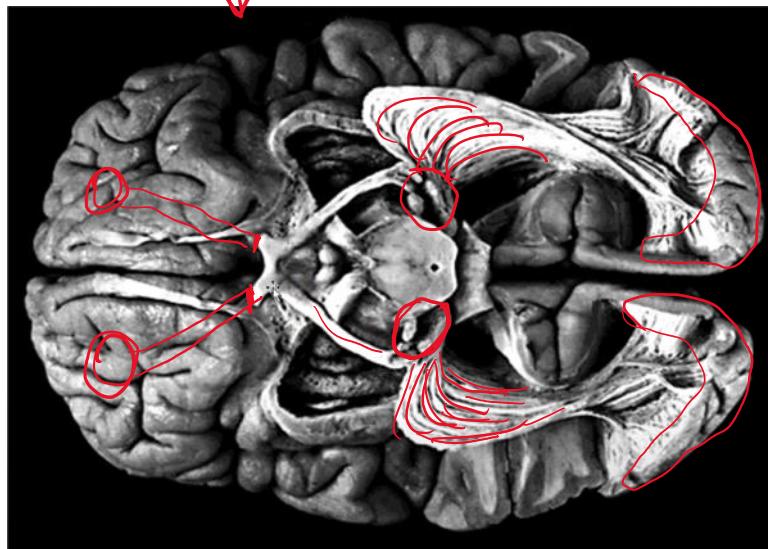
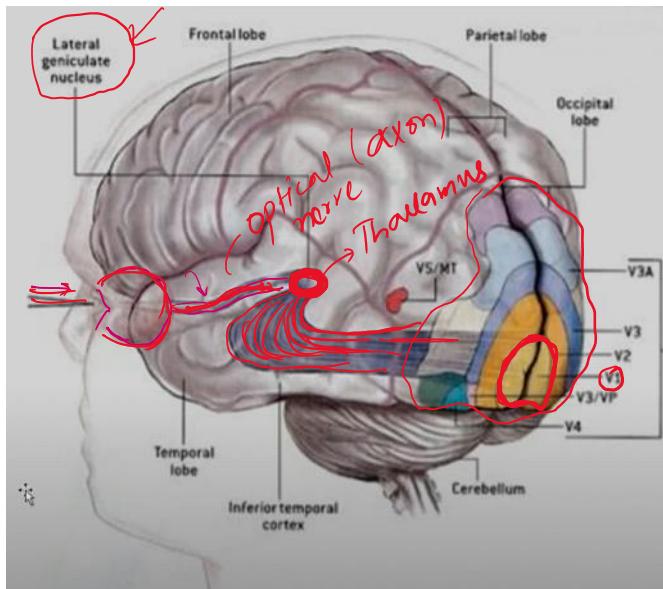


Human Visual Cortex

18 August 2022 16:05

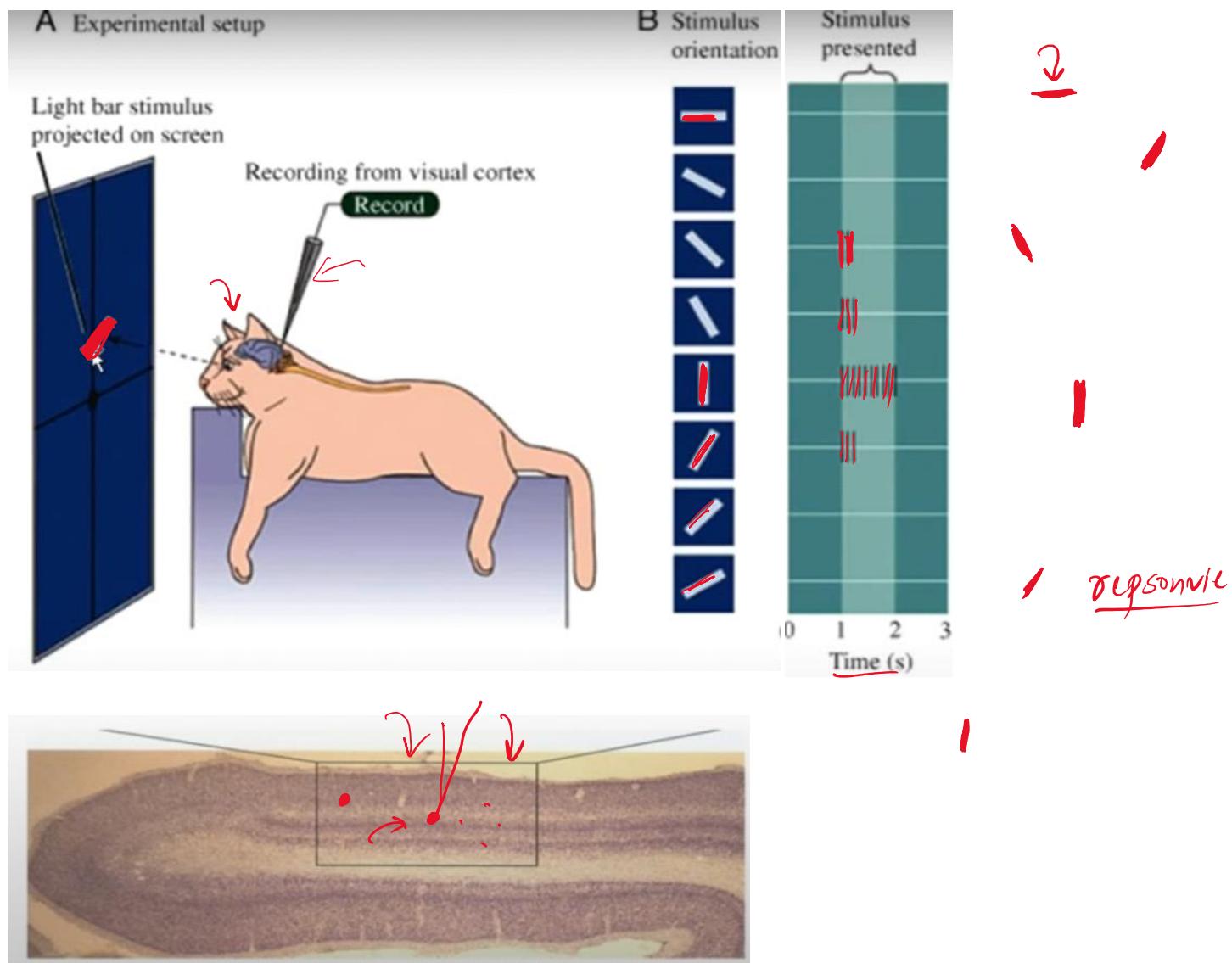
2D sheet primary V cortex

actual brain



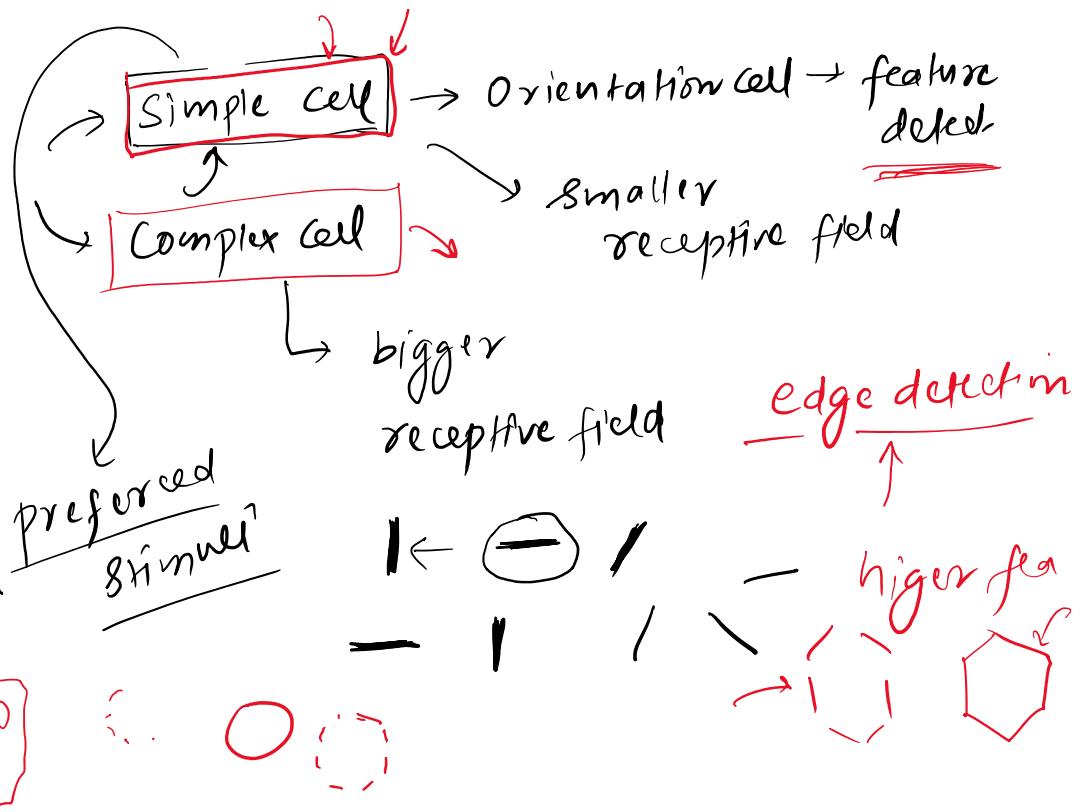
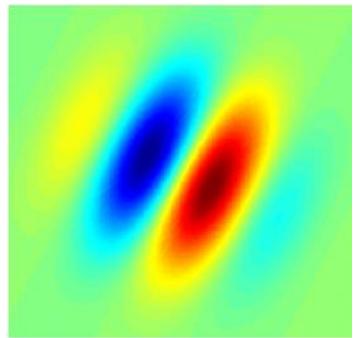
The Experiment

18 August 2022 16:09



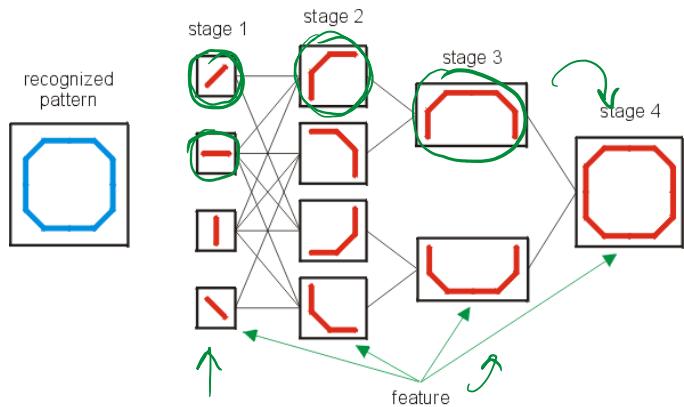
Conclusion

18 August 2022 17:57

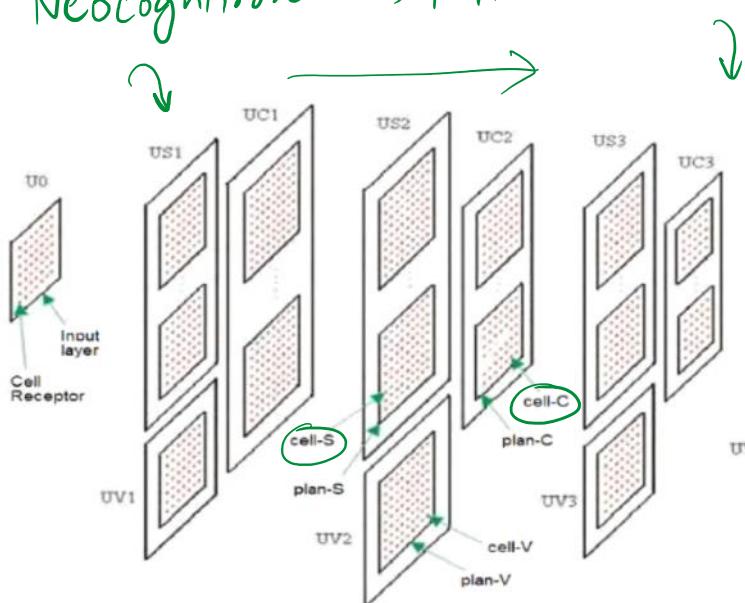


Development

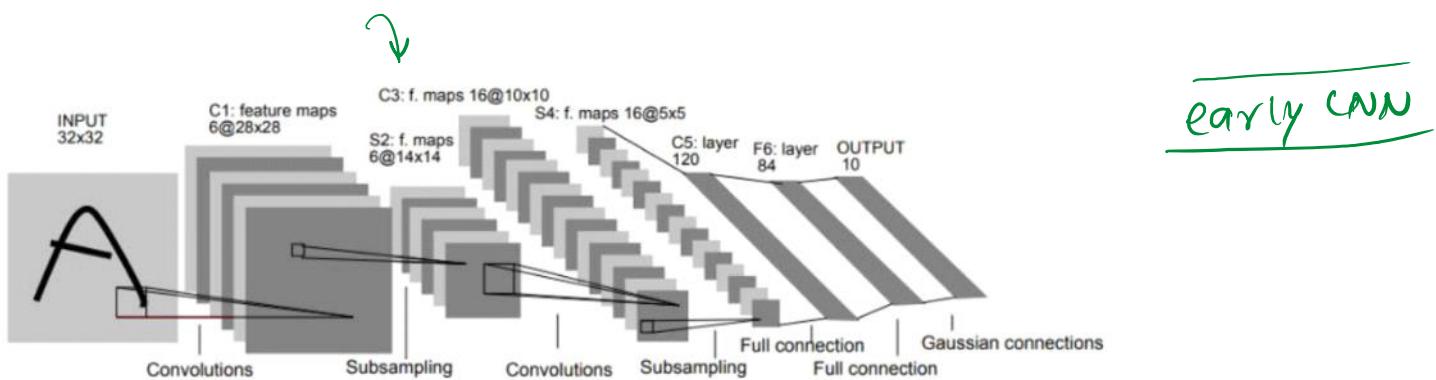
18 August 2022 16:15



Neocognitron → Fukushima



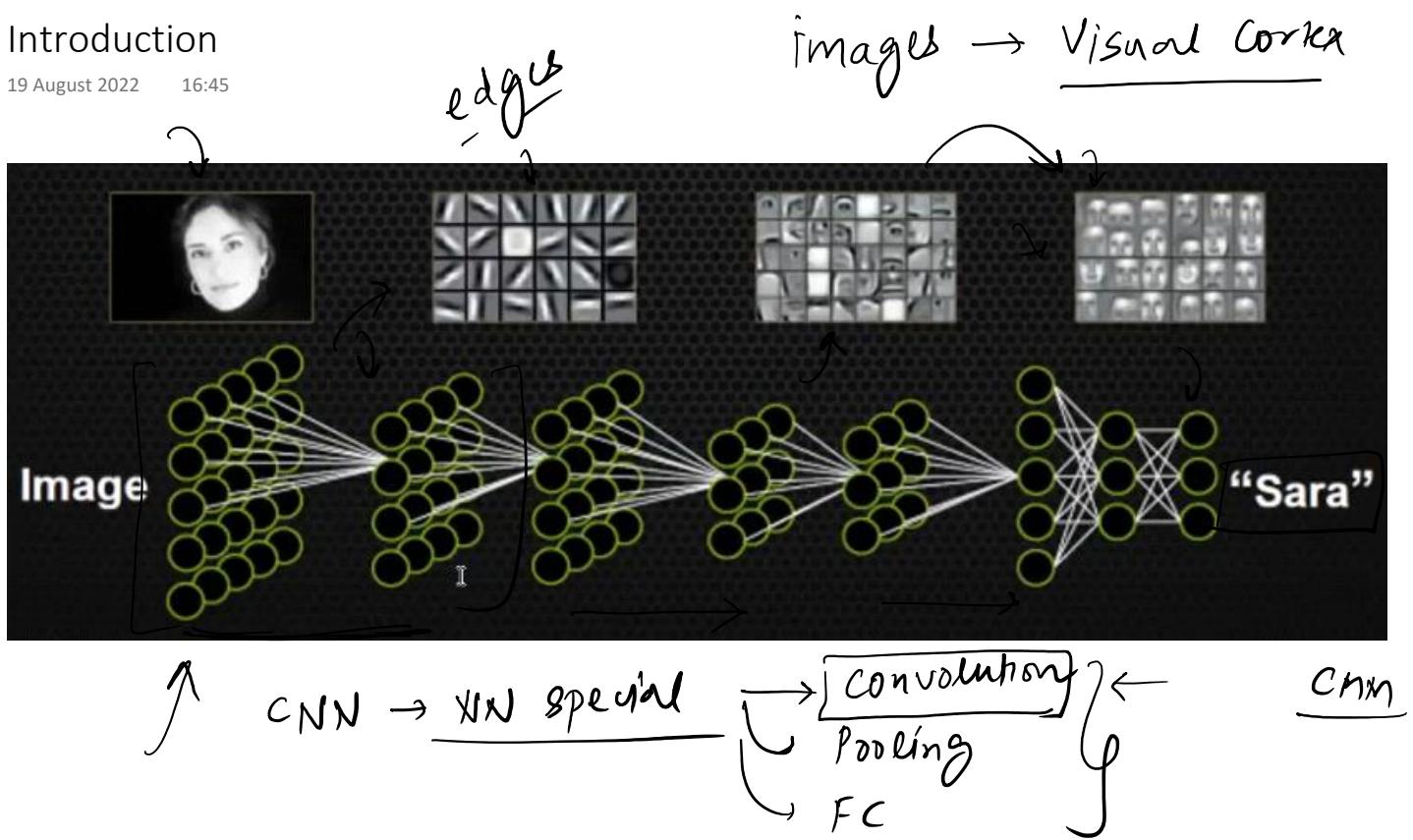
Yann LeCun → CNN → Backprop convolution



2012 → AlexNET → ImageNET → CNNs

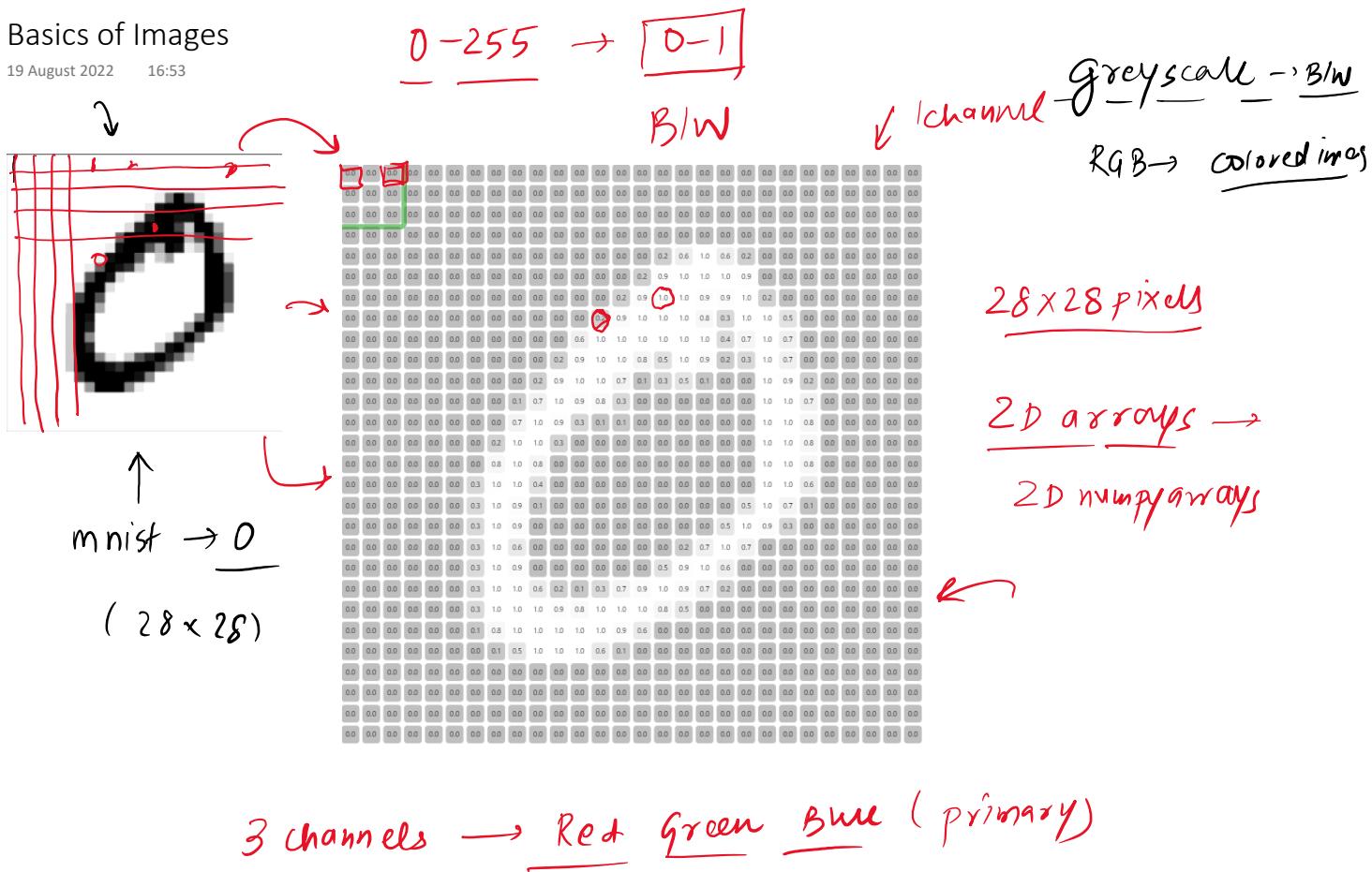
Introduction

19 August 2022 16:45

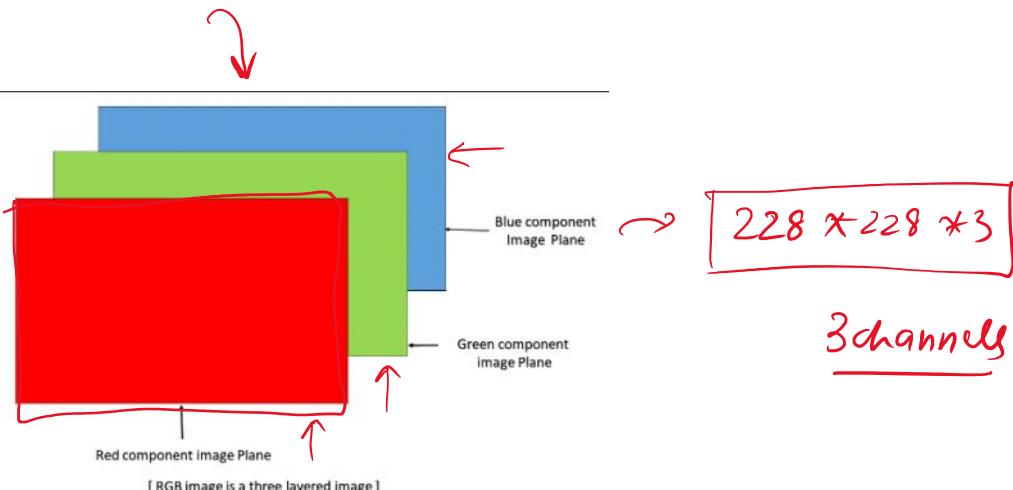


Basics of Images

19 August 2022 16:53

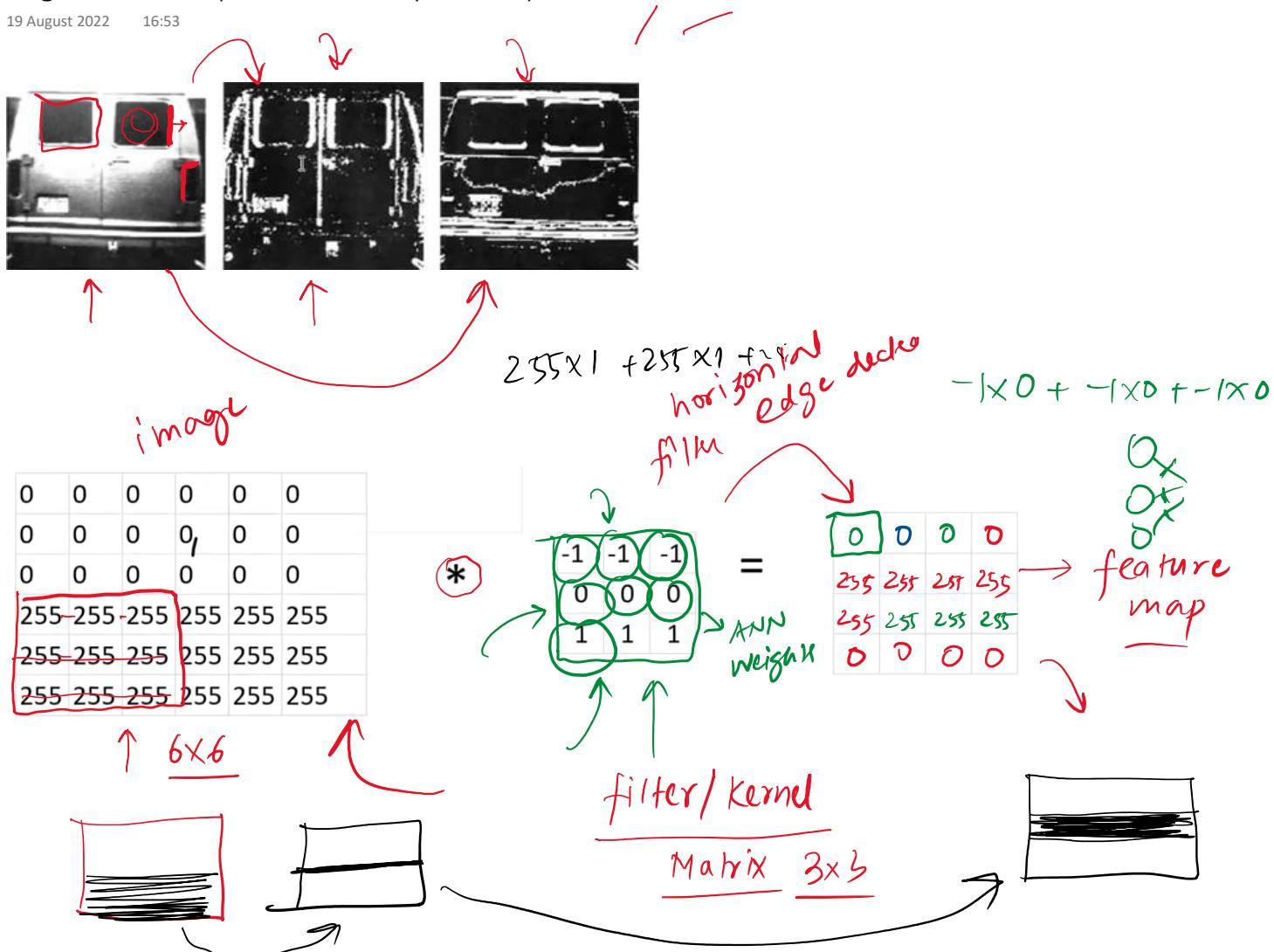


3 channels \rightarrow Red Green Blue (primary)



Edge Detection (Convolution Operation)

19 August 2022 16:53



$$\begin{array}{ccccccc}
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 0 & 0 & 0 & 0 & 0 & 0 \\
 255 & 255 & 255 & 255 & 255 & 255 \\
 255 & 255 & 255 & 255 & 255 & 255 \\
 255 & 255 & 255 & 255 & 255 & 255
 \end{array} * \begin{array}{ccc}
 -1 & -1 & -1 \\
 0 & 0 & 0 \\
 1 & 1 & 1
 \end{array} = \begin{array}{cccc}
 0 & 0 & 0 & 0 \\
 255 & 255 & 255 & 255 \\
 255 & 255 & 255 & 255 \\
 0 & 0 & 0 & 0
 \end{array}$$

(6×6) (3×3) (4×4)

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

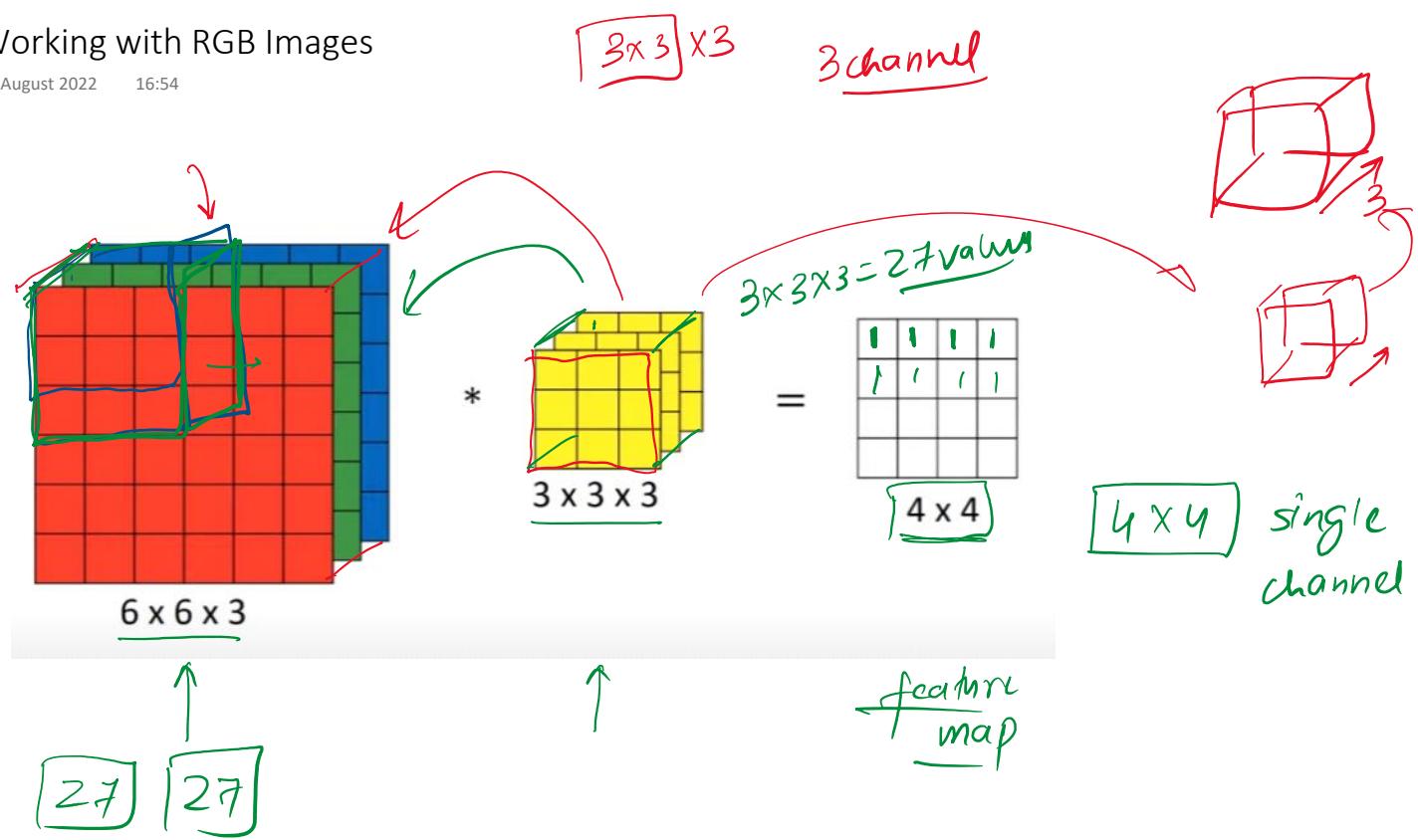
$$(64 \times 64) \quad (3 \times 3) \rightarrow \underline{(62 \times 62)}$$

Demo

19 August 2022 16:54

Working with RGB Images

19 August 2022 16:54



$$\boxed{m \times m \times c} \quad \boxed{n \times n \times c} \rightarrow \frac{(m-n+1) \times (m-n+1)}{\text{single channel}}$$

Multiple Filters

23 August 2022 08:24

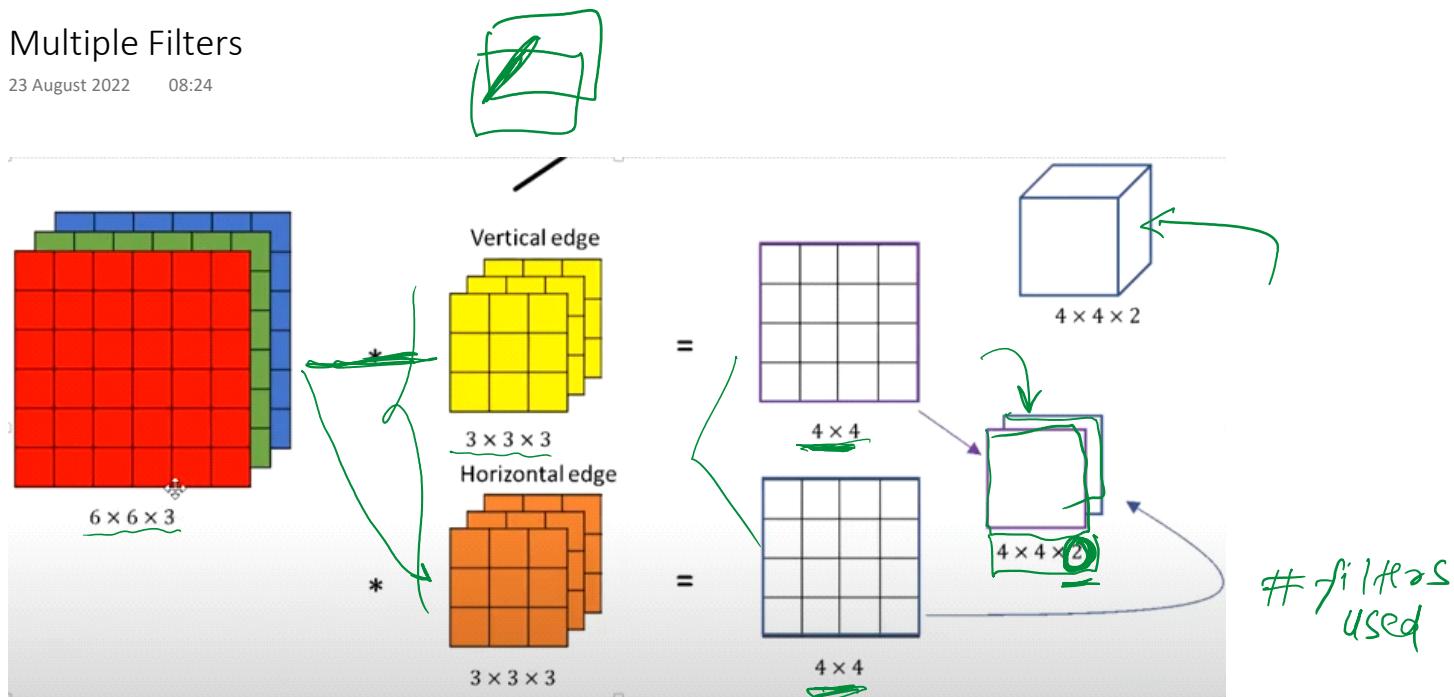


Image taken from Andrew NG's lecture

$$10 \rightarrow [4 \times 4 \times 10]$$

Problem with Convolution

26 August 2022 14:25

What is Padding?

26 August 2022 14:26

Handwritten annotations:

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

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

$$\frac{5 \times 5}{\text{---}} \rightarrow 3 \times 3$$

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

$(n - f + 1)$



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

$$5 + 2(1) - 3 + 1 = 7 - 3 + 1 = \underline{\circlearrowleft 5}$$

Keras

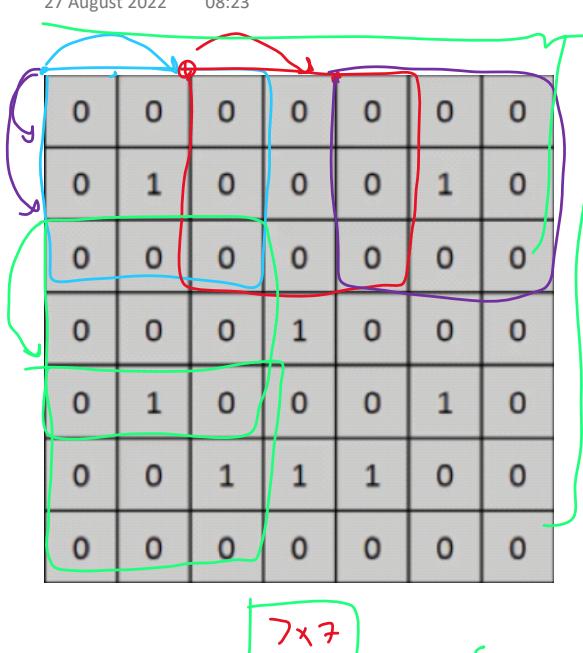
Valid

Same



Strides

27 August 2022 08:23



Stride = 1

0	0	1
1	0	0
0	1	1

3 x 3

$$\text{Stride} = (2, 2)$$

$$\text{Strichl} = 2 \rightarrow$$

$$\frac{t-5}{2} + 1$$

$$2+1=3$$

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

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

\downarrow

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

$$\begin{array}{r} \underline{7-3} \\ 2 \end{array} \qquad \begin{array}{r} \underline{6-3} \\ 2 \end{array} \qquad 1.5 = 1 + 1 = 2$$

Special Case

Stride = 2



2x3

$$\begin{array}{r} 6 \times 7 \\ \hline \end{array}$$

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

$$\begin{array}{c} 19 \rightarrow 1 \\ 1 \cdot) \rightarrow) \\ \downarrow \\ f \mid 00\gamma \end{array}$$

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

7x6

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

Why Strides are required?

27 August 2022 08:24

1) High level features

2) Computing →

Keras → stride

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

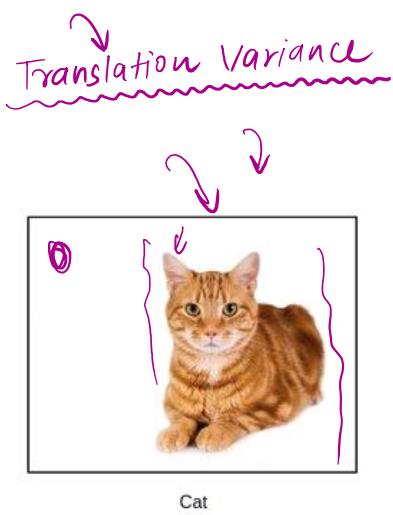
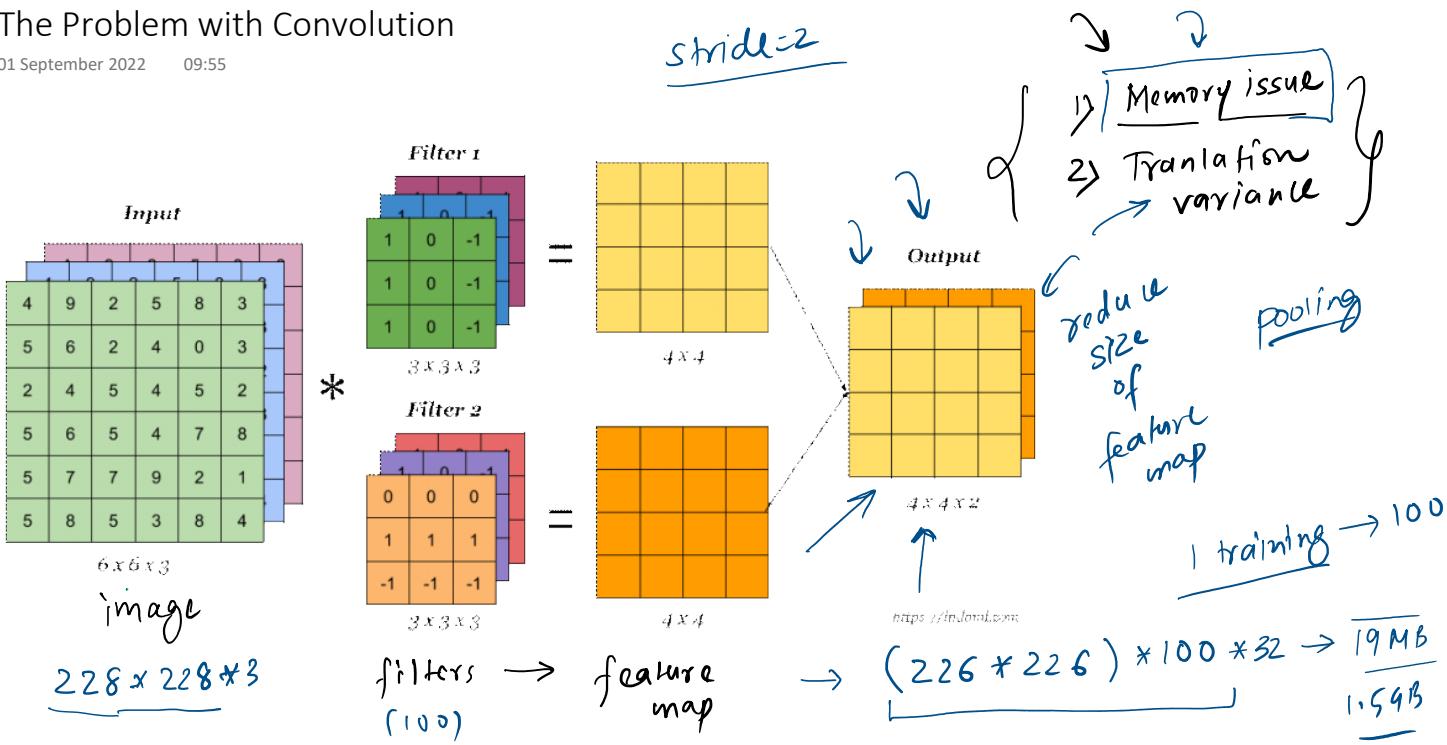
$$\frac{28 + 2 - 3}{2} + 1$$

$$\underline{13.5 + 1}$$

$$13 + 1 = 14$$

The Problem with Convolution

01 September 2022 09:55



The diagram illustrates the process of feature extraction and reduction:

- features** (underlined) are **location dependent**.
- location** (underlined) is indicated by an arrow pointing to the **features**.
- pooling** (underlined) is indicated by an arrow pointing to the **location**.
- A bracket labeled "down sample" covers the **pooling** and **your feature map** stages.

Pooling

01 September 2022 09:55

0	0	0	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	255	255	255

*

-1	-1	-1
0	0	0
1	1	1

\otimes_{clu}

(4x4)

Max pooling
Min pooling
Avg pooling
L2 pooling
Global pooling

-	-	-	-
-			

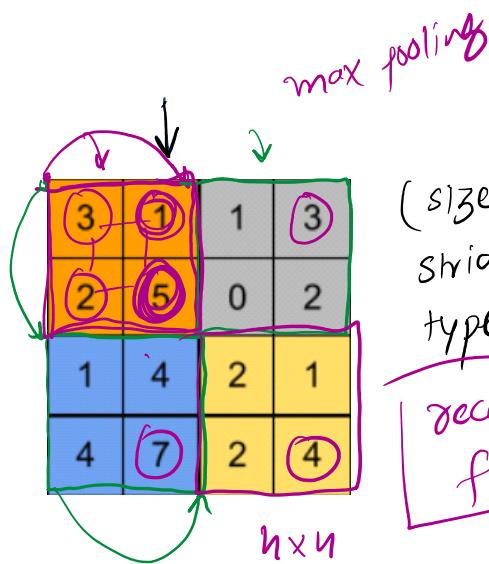
feature map (non-linear)

5	3
7	4

feature map

2x2

low level details eliminate



$\xrightarrow{\text{size}} (2,2)$
 $\xrightarrow{\text{stride}} (2)$
 $\xrightarrow{\text{type}} \text{Max}$

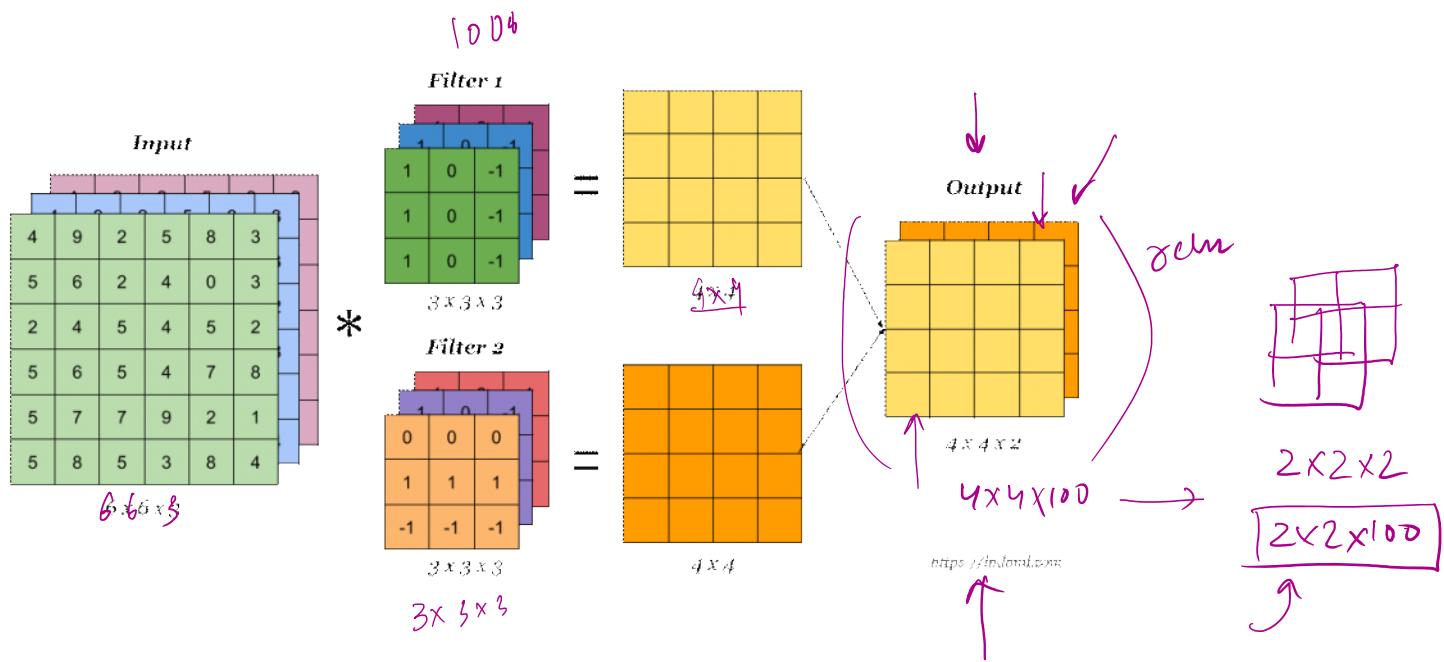
$\boxed{\text{receptive field}}$

Demo

01 September 2022 09:57

Pooling on Volumes

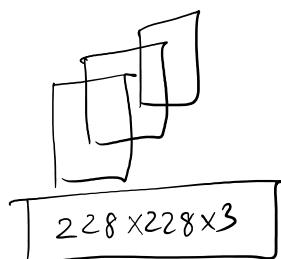
01 September 2022 09:56



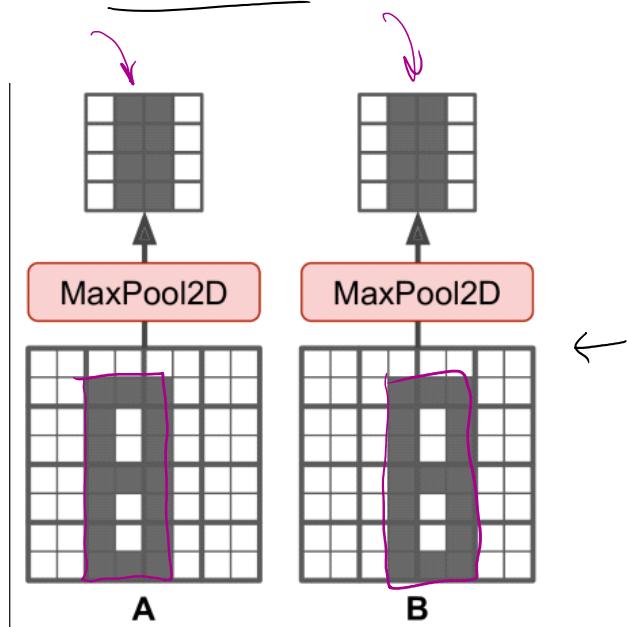
Advantages of Pooling

01 September 2022 09:56

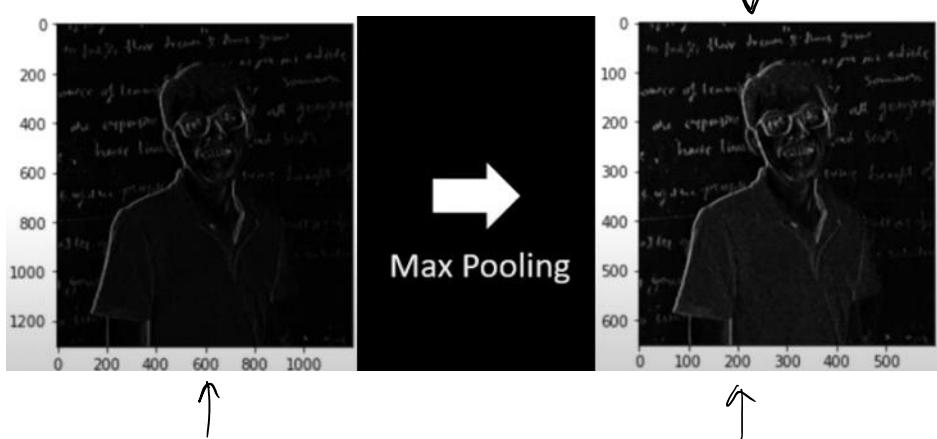
1) reduced size



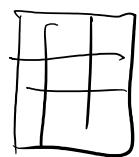
2) Translation invariance



3) Enhanced features
(only in case of Max pooling)



4) No need of training

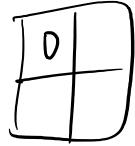


min

max pooling
avg pooling

faster

aggregate



(2×2)

2

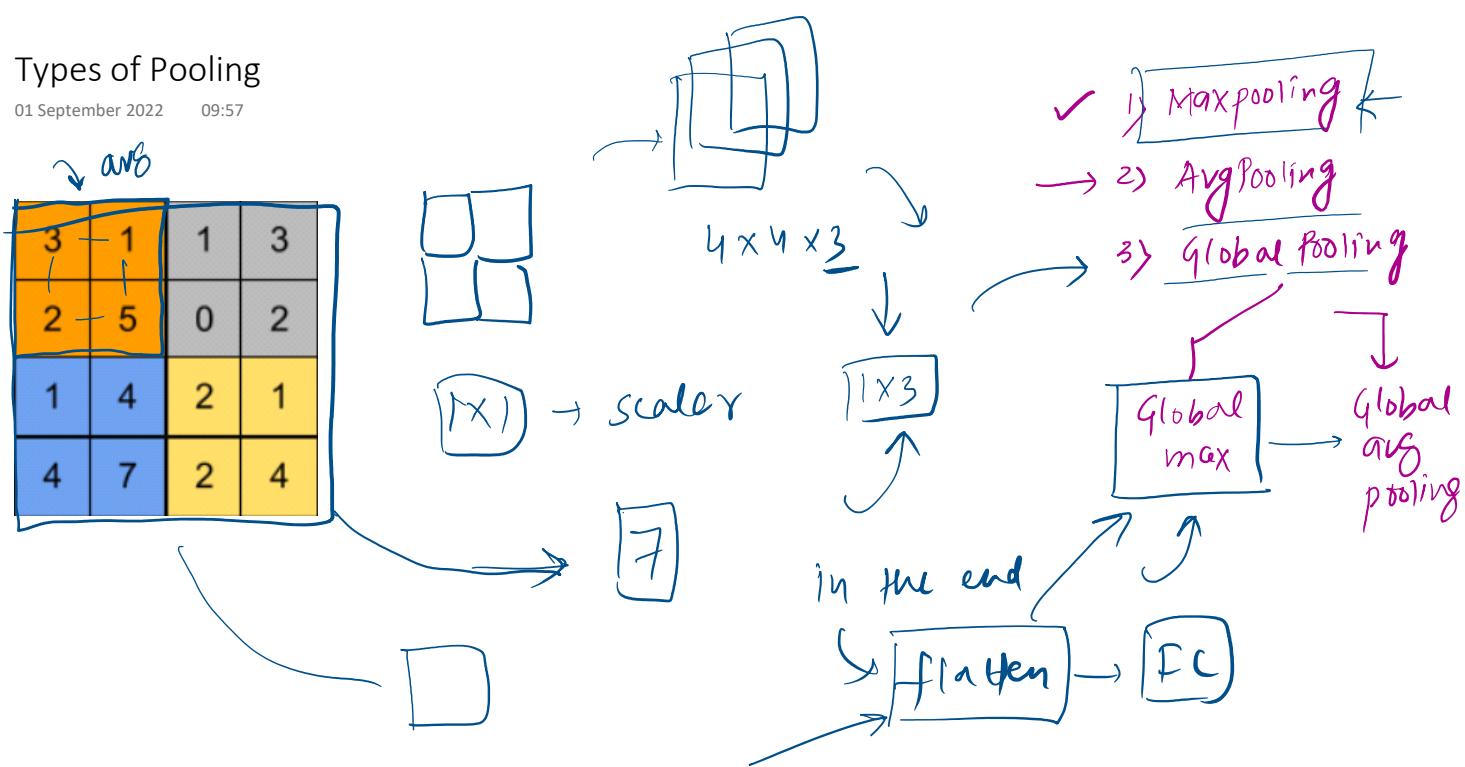
type

Keras Code

01 September 2022 09:56

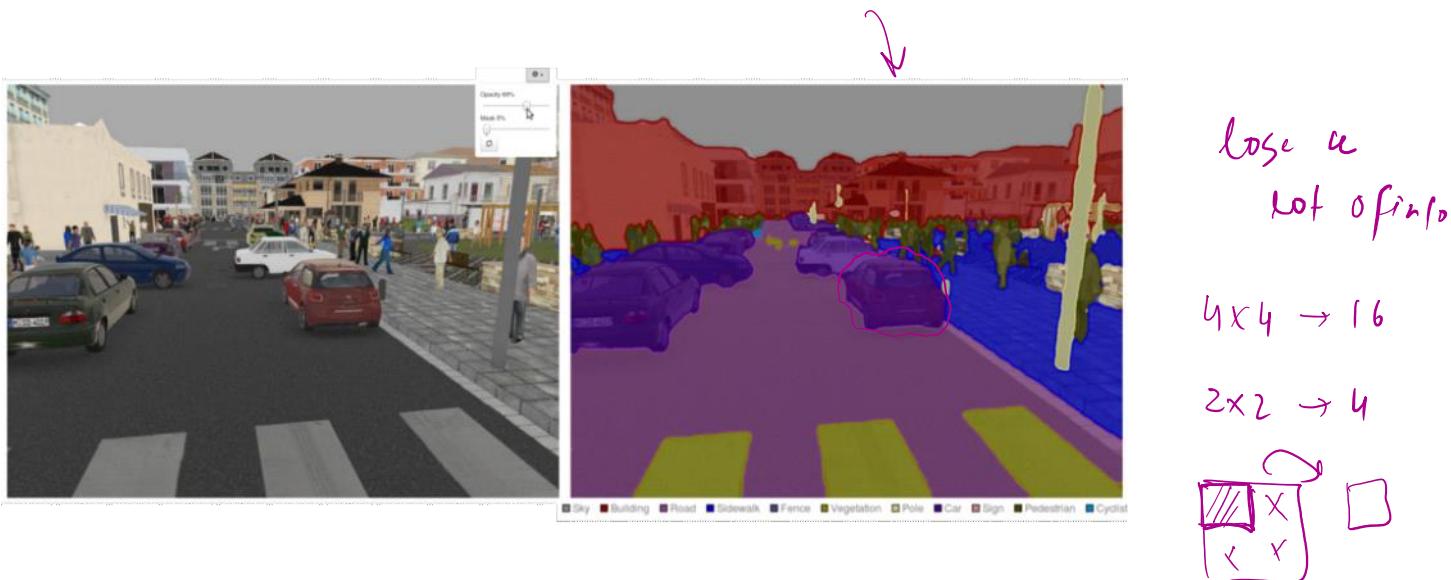
Types of Pooling

01 September 2022 09:57



Disadvantages of Pooling

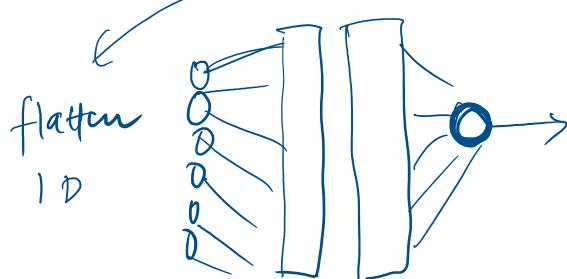
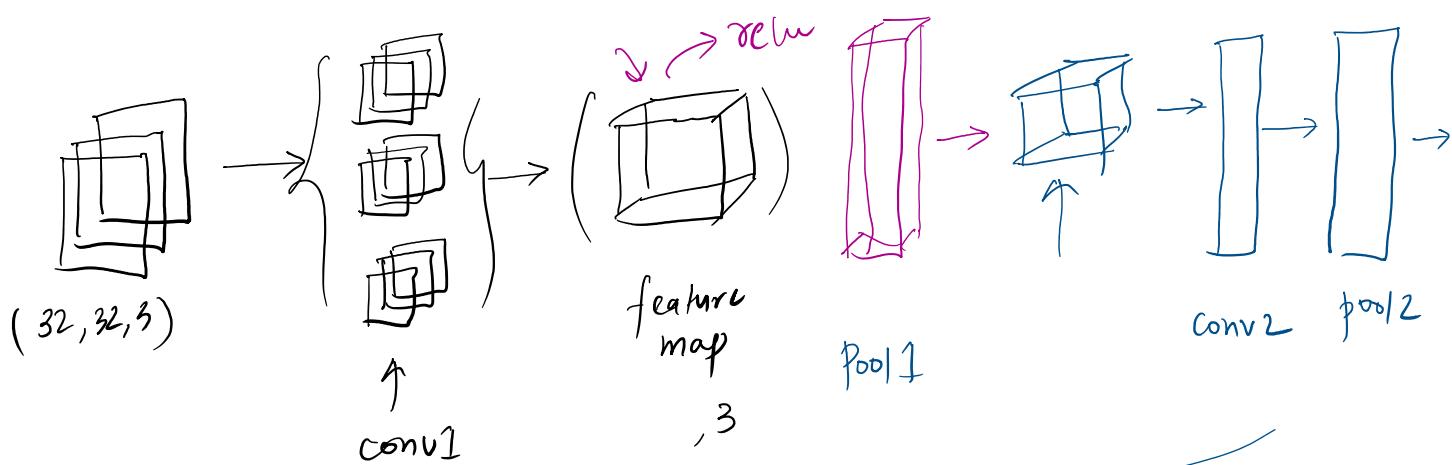
01 September 2022 09:57



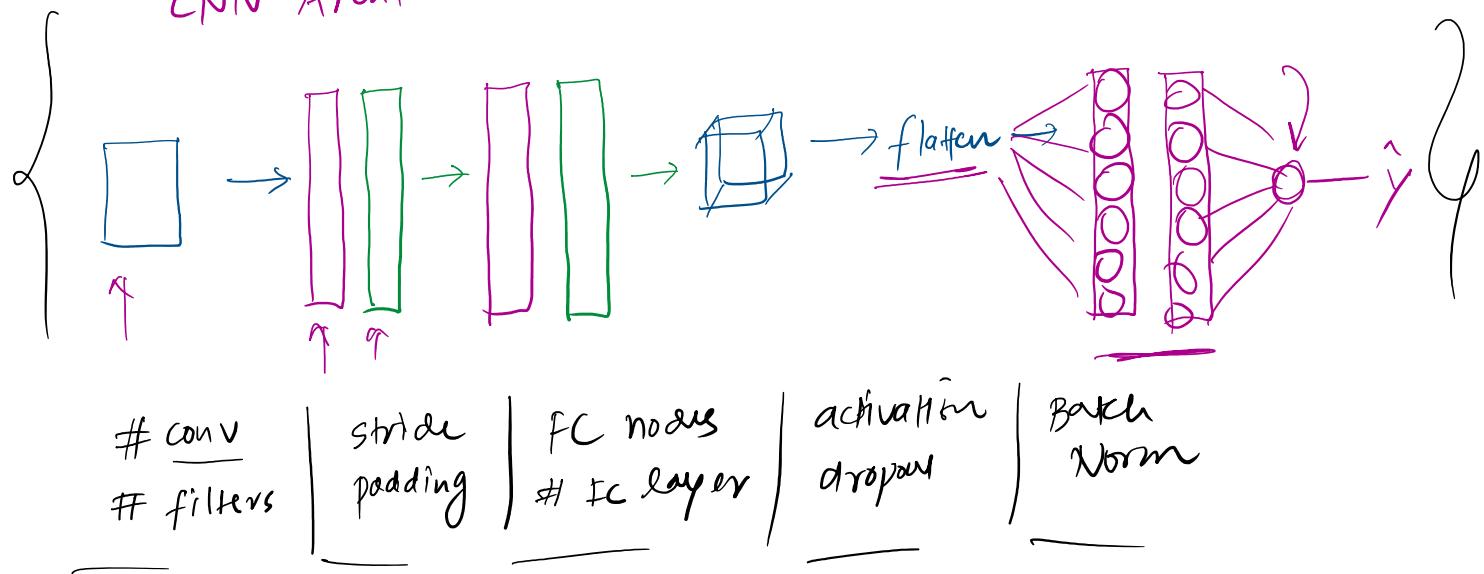
CNN Architecture

02 September 2022 10:55

1) Convolution 2) Padding / stride 3) Pooling



CNN Architecture



ImageNET

- 1) [LeNET] \rightarrow Yann LeCANN
- 2) AlexNET
- 3) GoogleLeNET

4) VggNET

5) ResNET

6) Inception

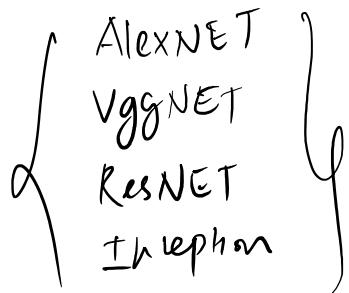
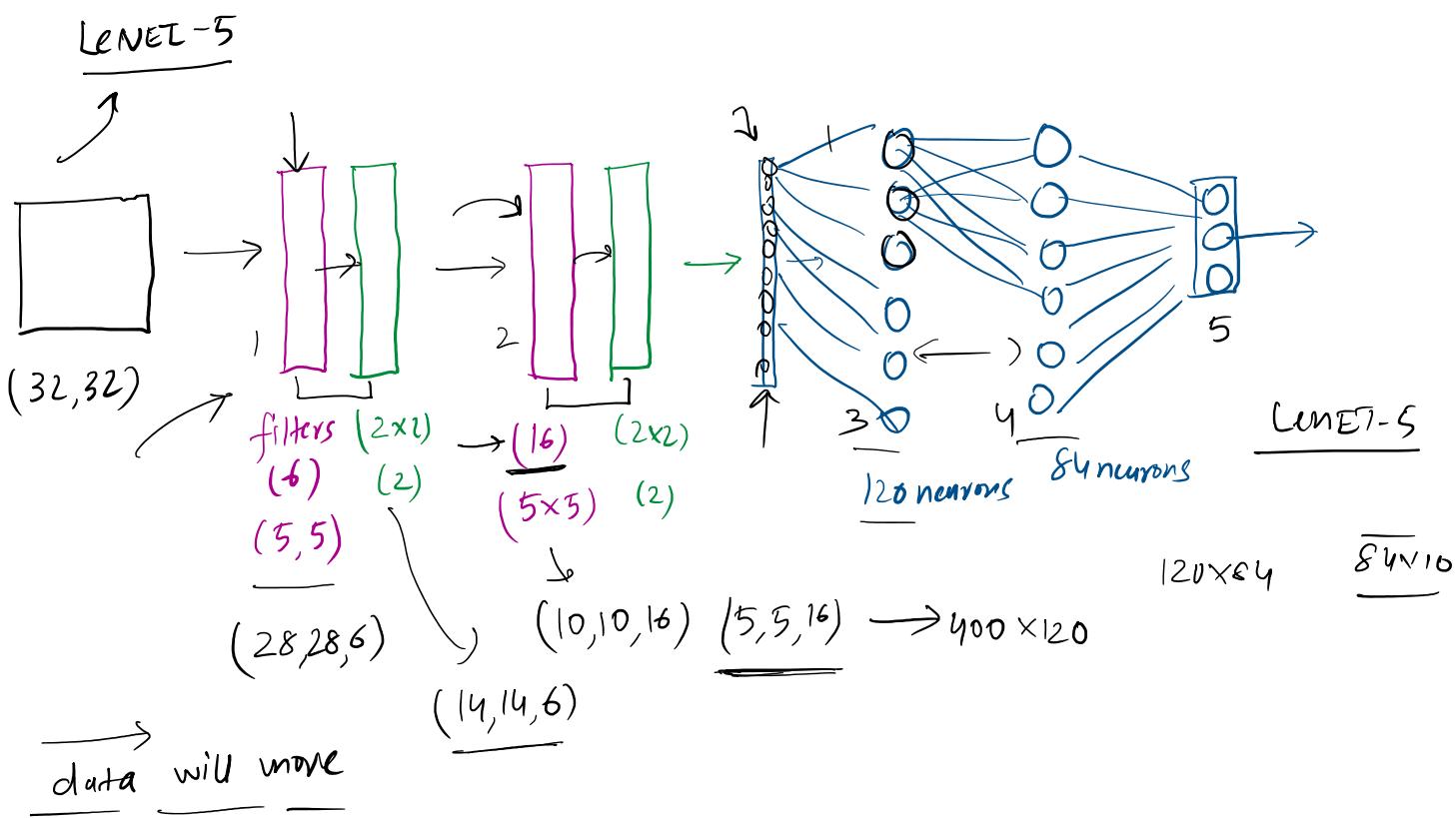
↳ GoogLeNET

↳ AlexNet

LeNet

02 September 2022 10:55

Yann LeCun → 1989 → 1998 → LeNET → CNN → US Navy Postal Service



Guidelines

02 September 2022 10:58

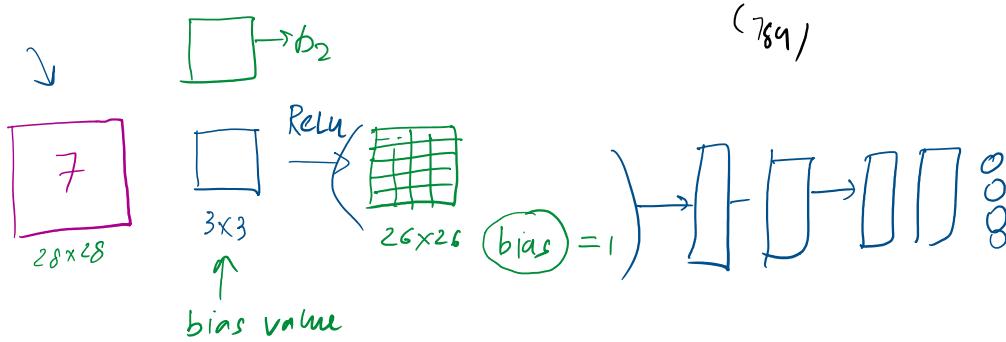
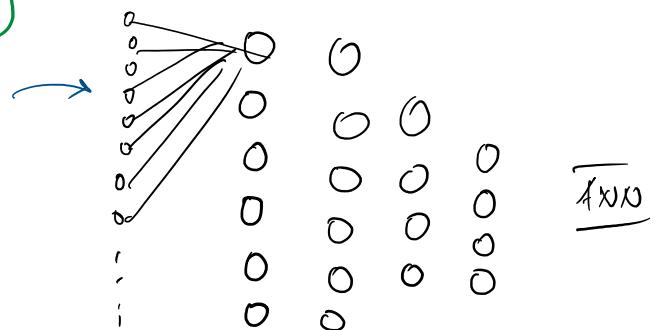
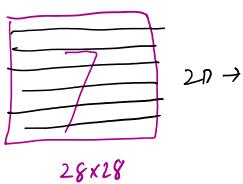
Keras Code

02 September 2022 14:58

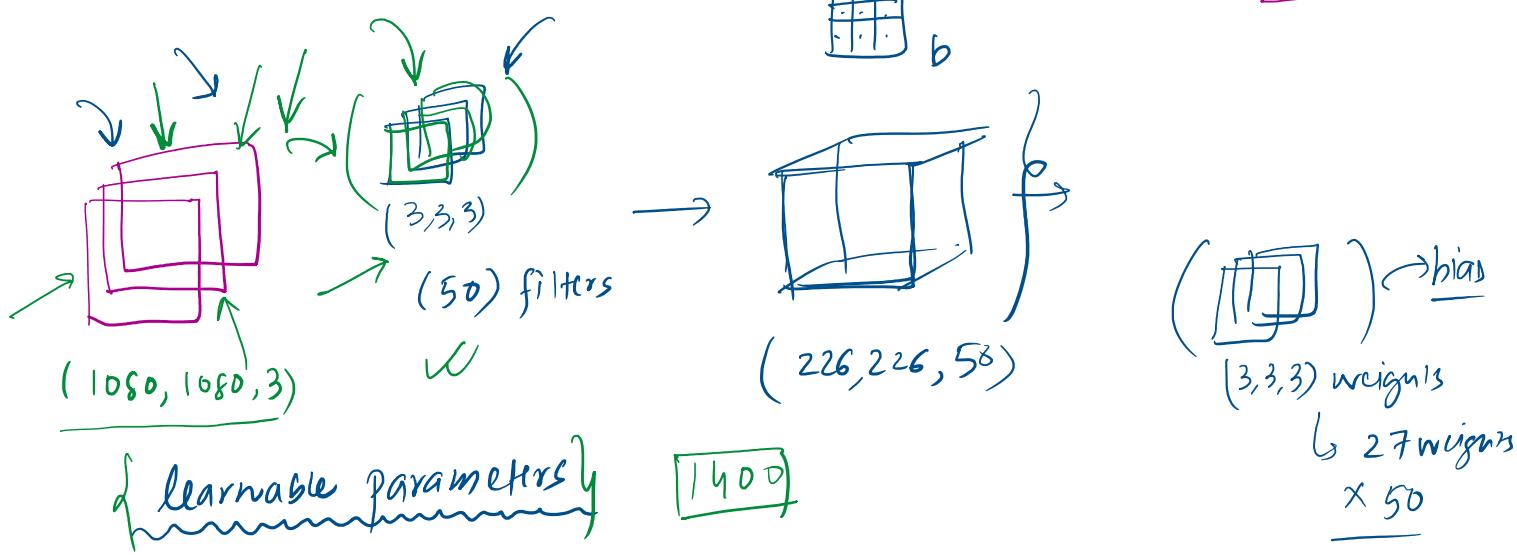
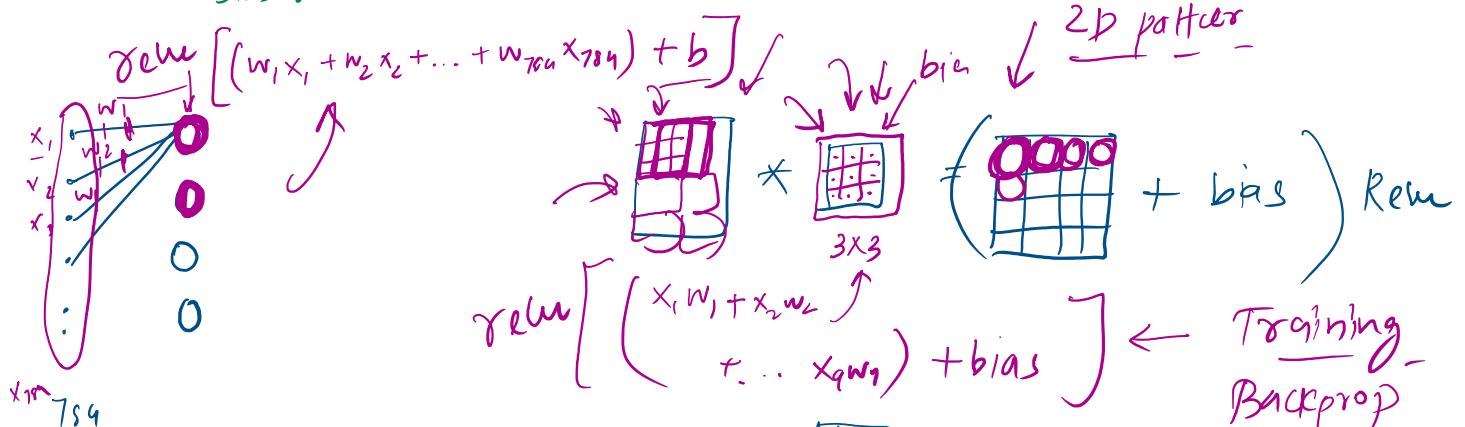
CNN Vs ANN

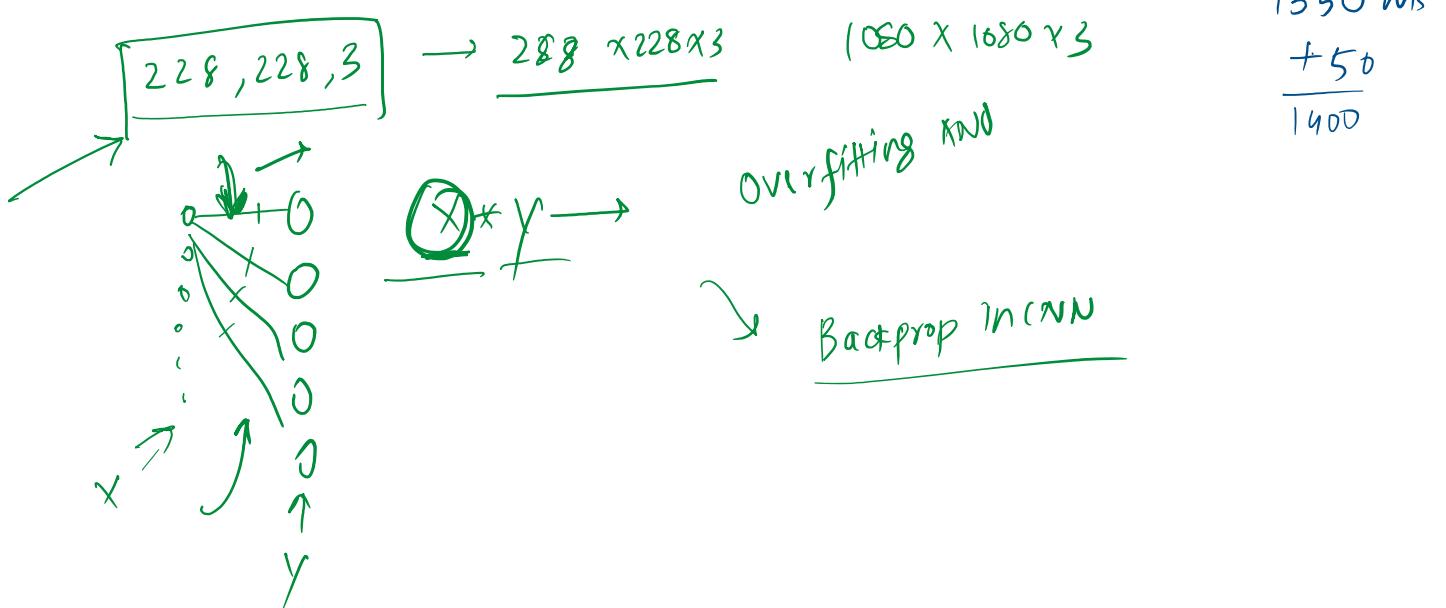
06 September 2022 10:00

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



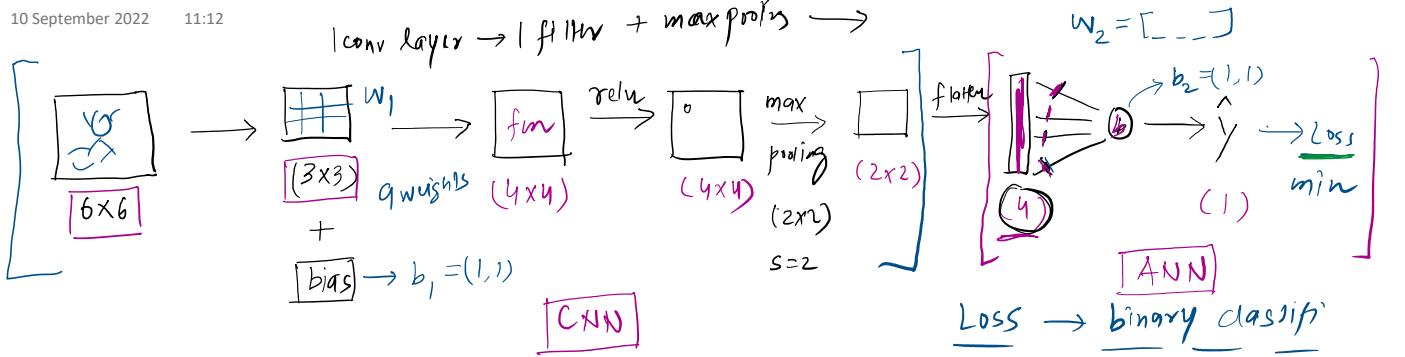
similarity ANN CNN





Backpropagation in CNN

10 September 2022 11:12



Trainable Parameters

$$W_1 = (3, 3) \quad W_2 = (1, 4) \quad = 15 \text{ trainable parameters}$$

$$b_1 = (1, 1) \quad b_2 = (1, 1)$$

Logical Flow

$$\text{Loss} \rightarrow \text{binary classif}$$

$$L = -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

$$A_2 = \hat{y}$$

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

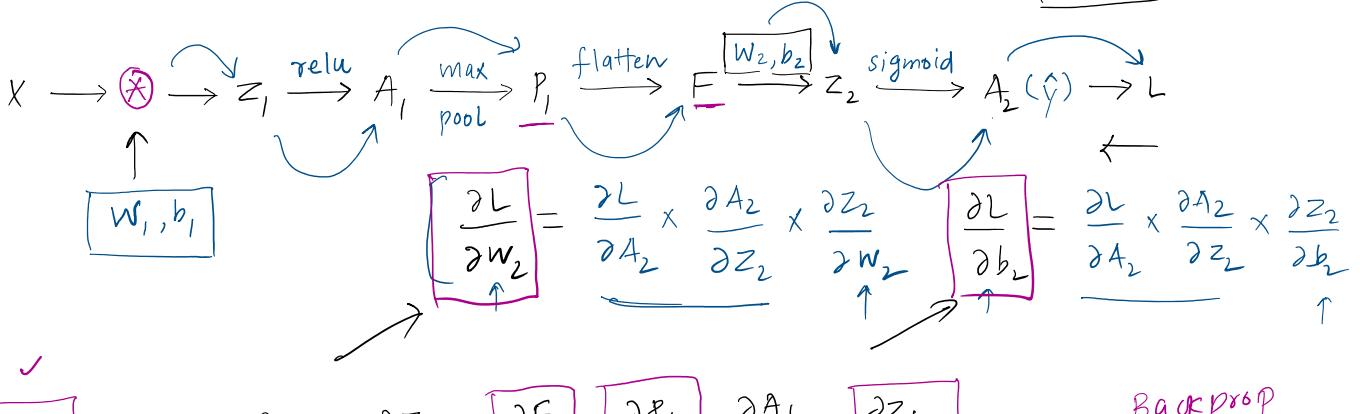
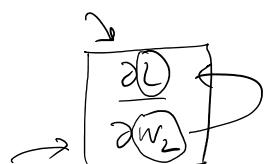
$$W_1 = W_1 - \eta \frac{\partial L}{\partial W_1}$$

$$W_2 = W_2 - \eta \frac{\partial L}{\partial W_2}$$

Loss is minimized

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

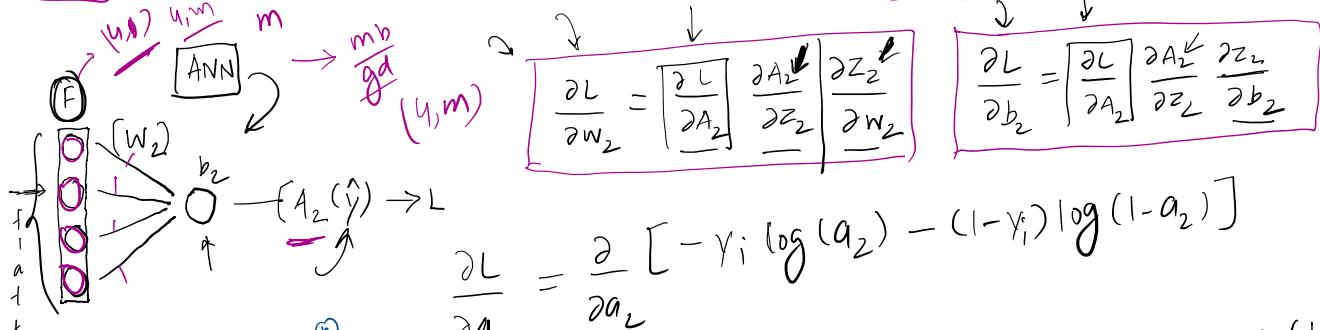
$$b_2 = b_2 - \eta \frac{\partial L}{\partial b_2}$$



$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial A_2} \times \frac{\partial A_2}{\partial Z_2} \times \frac{\partial Z_2}{\partial F} \times \left[\frac{\partial F}{\partial P_1} \right] \times \left[\frac{\partial P_1}{\partial A_1} \right] \times \frac{\partial A_1}{\partial Z_1} \times \left[\frac{\partial Z_1}{\partial w_1} \right]$$

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

Backprop
→ Convolution
→ Flatten
→ Max pooling



Forward Prop eqⁿ

$$\begin{cases} Z_2 = W_2 F + b_2 \\ A_2 = \sigma(Z_2) \end{cases}$$

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

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

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

w₂ update
shape =

$$\left[\frac{\partial Z_2}{\partial w_2} = F \right]$$

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

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

m images

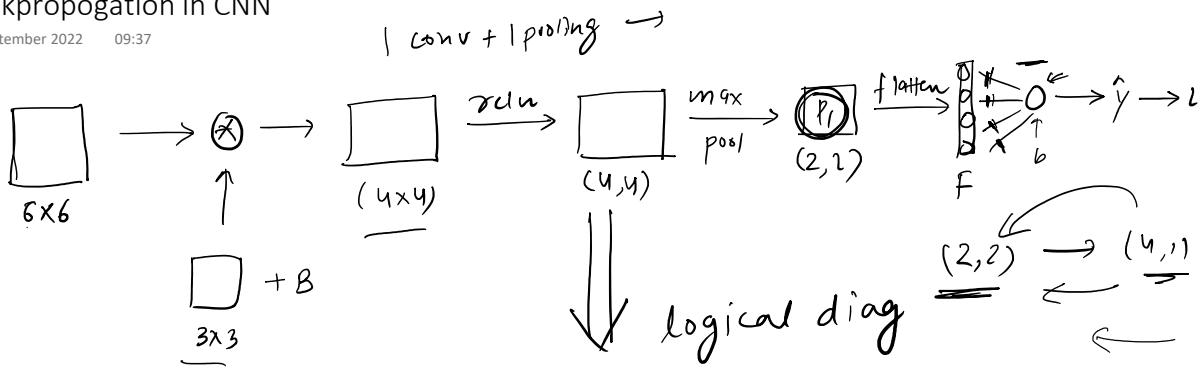
$$\left[\frac{\partial L}{\partial b_2} \right] = (A_2 - Y) F^T \quad \left| \quad \frac{\partial L}{\partial b_2} = (A_2 - Y) \right.$$

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

\uparrow $(1, m) - (1, m)$
 \downarrow
 $(1, m) \quad (m, n) \rightarrow \boxed{(1, 4)} \rightarrow \underbrace{w_2}_{\text{batch of image}} \rightarrow (1, n)$

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

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

$$\left[\frac{\partial z_2}{\partial F} \right] = w_2 \rightarrow$$

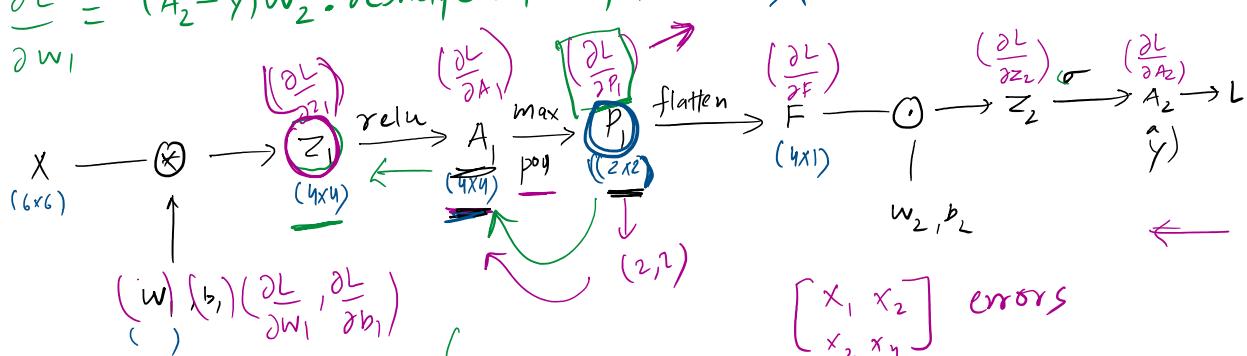
Shape? $\rightarrow (F)$

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

reshape(P1.shape)

$$\frac{\partial L}{\partial w_1} = (A_2 - y) w_2 \cdot \text{reshape}(P_1, \text{shape})$$

$\frac{\partial L}{\partial w_1}$

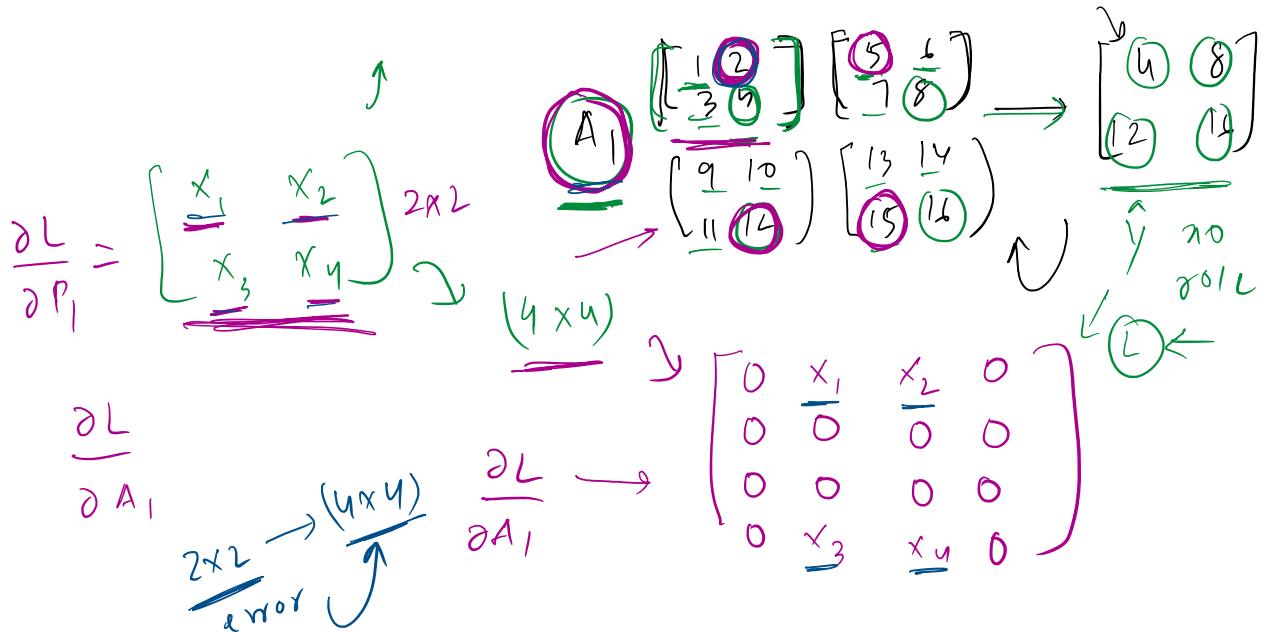


$\begin{bmatrix} x_1 & x_2 \\ x_3 & x_4 \end{bmatrix}$ errors

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

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

flatten \rightarrow no trainable parameters



$\frac{\partial L}{\partial w_2} = \left[\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial z_2} \frac{\partial z_2}{\partial F} \frac{\partial F}{\partial p_1} \frac{\partial p_1}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \right]$

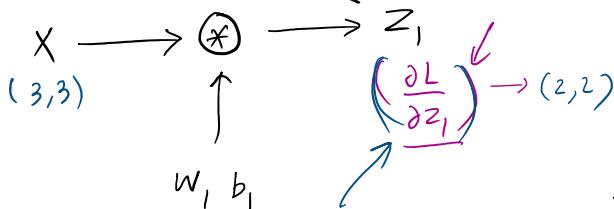
$\frac{\partial L}{\partial b_1} = \left[\frac{\partial L}{\partial A_2} \frac{\partial A_2}{\partial z_2} \frac{\partial z_2}{\partial F} \frac{\partial F}{\partial p_1} \frac{\partial p_1}{\partial A_1} \frac{\partial A_1}{\partial z_1} \frac{\partial z_1}{\partial b_1} \right]$

$\frac{\partial L}{\partial A_1} = \left\{ \begin{array}{ll} \frac{\partial L}{\partial p_1}_{xy}, & \text{if } A_{xy} \text{ is the max element} \\ 0, & \text{otherwise} \end{array} \right.$

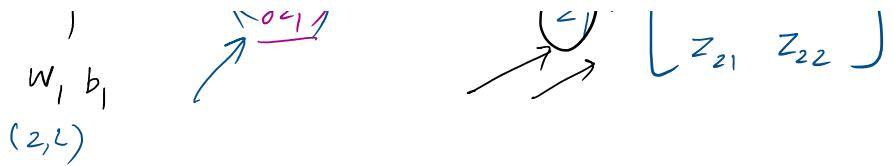
$\frac{\partial A_1}{\partial z_1} = \left\{ \begin{array}{ll} 1 & \text{if } z_{1xy} > 0 \\ 0 & \text{if } z_{1xy} < 0 \end{array} \right.$

Convolution $\xrightarrow{\text{Backprop}}$ max pooling $\xrightarrow{\text{fatten}}$

Backprop on Convolution



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



$$X = \underbrace{\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{bmatrix}}_{\text{---}} \times \underbrace{\begin{bmatrix} x_{13} \\ x_{23} \\ x_{33} \end{bmatrix}}_{\text{---}} \otimes \underbrace{W_1 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}}_{\text{---}} + \underbrace{b_1}_{\text{---}}$$

$$\boxed{\frac{\partial L}{\partial b_1}} = \frac{\partial L}{\partial z_1} \times \frac{\partial z_1}{\partial b_1} = \left(\frac{\partial L}{\partial z_{11}} \frac{\partial z_{11}}{\partial b_1} + \frac{\partial L}{\partial z_{12}} \frac{\partial z_{12}}{\partial b_1} + \frac{\partial L}{\partial z_{21}} \frac{\partial z_{21}}{\partial b_1} + \frac{\partial L}{\partial z_{22}} \frac{\partial z_{22}}{\partial b_1} \right)$$

↑ ↑ ↑ ↑

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

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

bias

$$X \rightarrow \otimes \rightarrow \boxed{Z_1} \quad \left(\frac{\partial L}{\partial z_1} \right)$$

↑ ↗

(3×3) W_1, b_1 (2×2)

$$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix} \quad W_1 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} + \dots$$

$$\frac{\partial L}{\partial W_1} = \begin{bmatrix} \frac{\partial L}{\partial w_{11}} & \frac{\partial L}{\partial w_{12}} \\ \frac{\partial L}{\partial w_{21}} & \frac{\partial L}{\partial w_{22}} \end{bmatrix} \quad \frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$$

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

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

↓

$$\frac{\partial L}{\partial w_{21}} = \underbrace{\frac{\partial L}{\partial z_{11}} \times \frac{\partial z_{11}}{\partial w_{21}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} \times \frac{\partial z_{12}}{\partial w_{21}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} \times \frac{\partial z_{21}}{\partial w_{21}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} \times \frac{\partial z_{22}}{\partial w_{21}}}_{\text{green bracket}}$$

$$\frac{\partial L}{\partial w_{22}} = \underbrace{\frac{\partial L}{\partial z_{11}} \times \frac{\partial z_{11}}{\partial w_{22}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} \times \frac{\partial z_{12}}{\partial w_{22}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} \times \frac{\partial z_{21}}{\partial w_{22}}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} \times \frac{\partial z_{22}}{\partial w_{22}}}_{\text{green bracket}}$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial w_{11}} = \underbrace{\frac{\partial L}{\partial z_{11}} x_{11}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} x_{12}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} x_{21}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} x_{22}}_{\text{green bracket}} \\ \frac{\partial L}{\partial w_{12}} = \underbrace{\frac{\partial L}{\partial z_{11}} x_{12}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} x_{13}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} x_{22}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} x_{23}}_{\text{green bracket}} \\ \frac{\partial L}{\partial w_{21}} = \underbrace{\frac{\partial L}{\partial z_{11}} x_{21}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} x_{22}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} x_{31}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} x_{32}}_{\text{green bracket}} \\ \frac{\partial L}{\partial w_{22}} = \underbrace{\frac{\partial L}{\partial z_{11}} x_{22}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{12}} x_{23}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{21}} x_{32}}_{\text{green bracket}} + \underbrace{\frac{\partial L}{\partial z_{22}} x_{33}}_{\text{green bracket}} \end{array} \right\}$$

$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$ $\frac{\partial L}{\partial z_1} = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} & \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{21}} & \frac{\partial L}{\partial z_{22}} \end{bmatrix}$

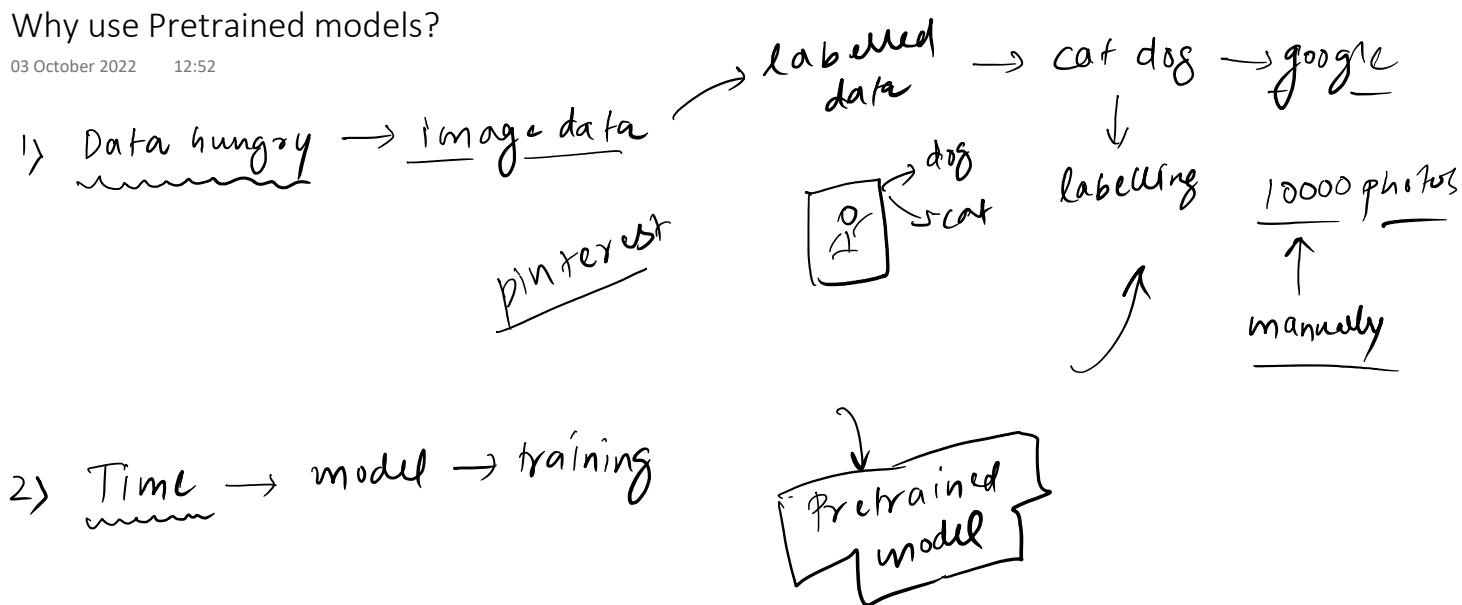
$\frac{\partial L}{\partial w_1} = \text{conv}(X, \frac{\partial L}{\partial z_1})$

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

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

Why use Pretrained models?

03 October 2022 12:52

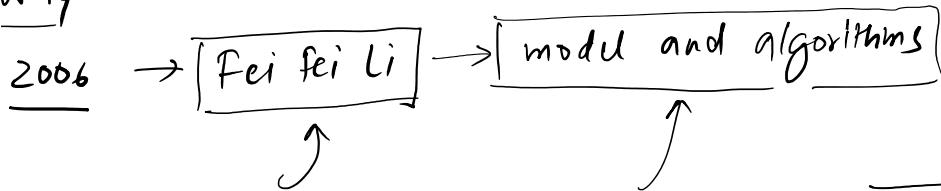


ImageNET Dataset

03 October 2022 12:36

Visual Database of images (Why What and How)

Why

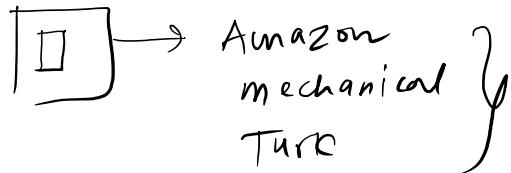


1.4 million images → 1.4 crore image → 20,000 categories → labels

1 million images → to batch → bounding box
↓ object local

```
graph LR; A[1.4 million images] --> B[1.4 crore images]; B --> C[20,000 categories]; C --> D[labels]; E[1 million images] --> F[batch]; F --> G[bounding box]; G --> H[object local];
```

How → crowd sourcing



→ Dataset → Deep learning

ImageNET challenge

