# NPTEL Week 10 Live Sessions

## on Deep Learning (noc24_ee04)

**A course offered by: Prof. Prabir Kumar Biswas, IIT Kharagpur**

- **Quiz 9 Solution**

Week 10 → Min-Max
         → Z-score
         → Batch Norm
         → Layer, Instance, group

→ Optimizers → Adam
            → Adgrad
            → RMS Prop
            → Momentum
            → NAG

**By**

**Arka Roy**
**NPTEL PMRF TA**

Prime Minister's Research Fellow
Department of Electrical Engineering, IIT Patna
Web: https://sites.google.com/view/arka-roy/home

① Comment on the learning rate of (Adagrad). Choose the correct option.

$grad_v \gg grad_h$.

Adaptive gradient optimizer.

a. Learning rate is adaptive

b. Learning rate increases for each time step

c. Learning rate remains the same for each update

d. None of the above

$L(w)$

$L(w_1, w_2)$

$grad_h$

$\downarrow r_t \uparrow$

$\uparrow \frac{1}{r_t}$

$\downarrow \downarrow$

Top view

$\downarrow \uparrow \eta = \frac{\eta}{\sqrt{\epsilon + r_t I}}$ Multiply

$g_t \downarrow$
$r_t \uparrow$
$\frac{1}{r_t} \downarrow$

$\frac{1}{r_t} \uparrow$ element wise

$w_2$

$grad_h$

$L(w_1, w_2)$

$r_2$

Equal energy line

$r_t = \sum g_t \circ g_t$

$grad_v$

$\Theta$

$r_1$

$g_t = \frac{1}{N} \sum_{x \in N \, minibatch} \nabla L(w)$

$r_t = g_t^2$

$W_{n+1} = w_n - \eta \nabla_w L(w) + \sum V_{t-1}$

Momentum optimizer.

Avg.

$W_{n+1} = w_n - \frac{\eta}{\sqrt{\epsilon + r_t I}} g_t$

# For the following figure A and figure B of loss landscape, choose correct statement



Figure A

Figure B

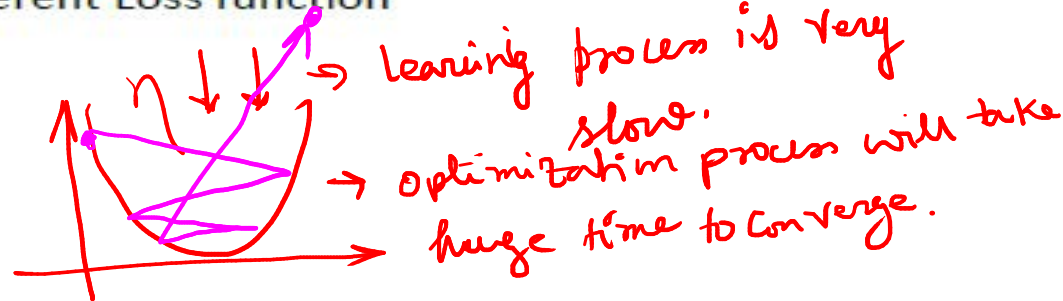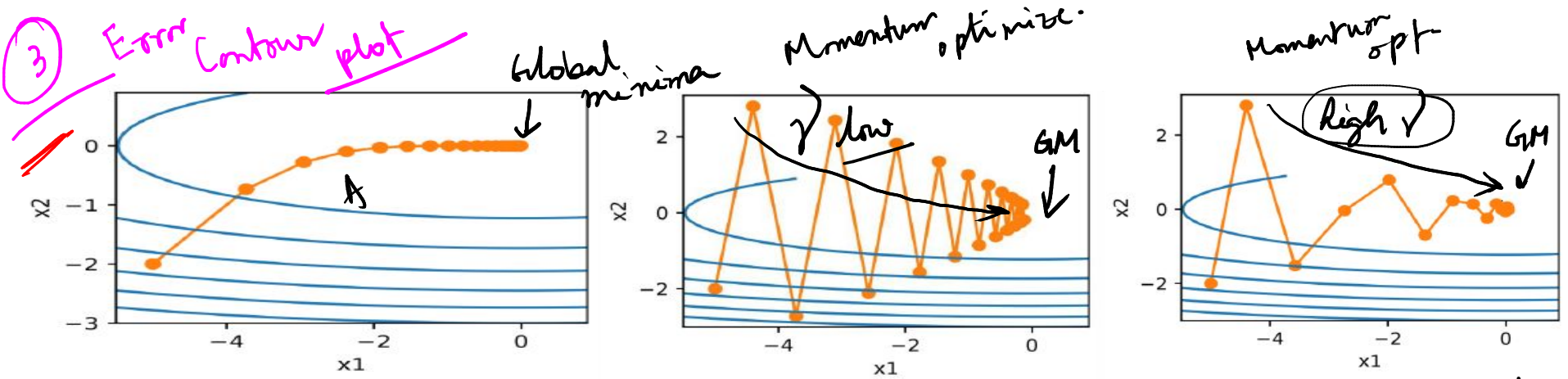a) Figure A has very small learning rate, Figure B has optimal learning rate
b) Figure A has optimal learning rate, Figure B has very small learning rate
c) Figure A and Figure B have different Loss function
d) None of Above

$$W_{n+1} = W_n - \eta \nabla_w L(w)$$

→ Learning rate

Loss/feature

Gradient descent rate/optimization Technique

highly oscillating around the optimal/global minima

diverging

$\eta \to$ high

$\eta \to$ low

Softmax

$\eta \downarrow \downarrow$ → learning process is very slow.
→ optimization process will take huge time to converge.

Optimal

③ **Error Contour plot**

Global minima
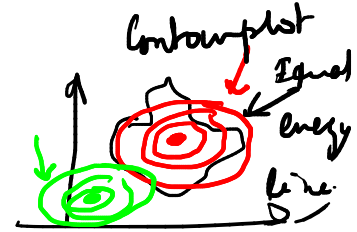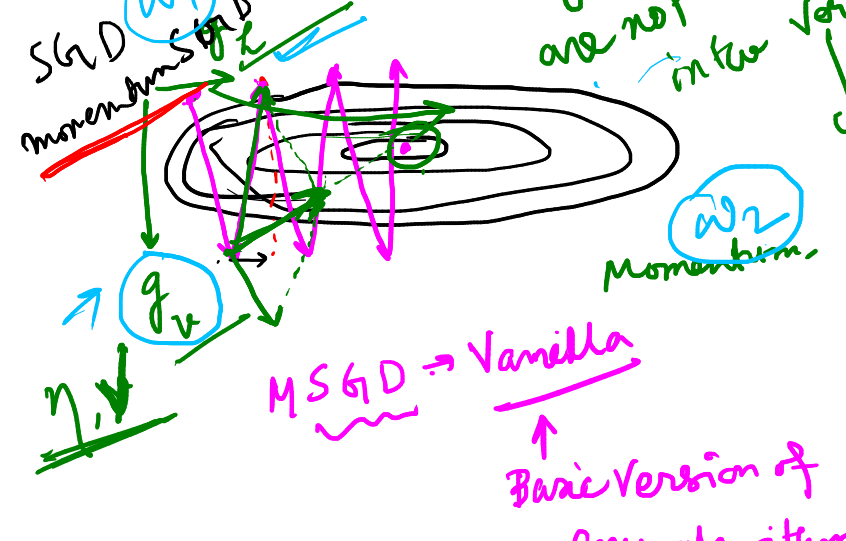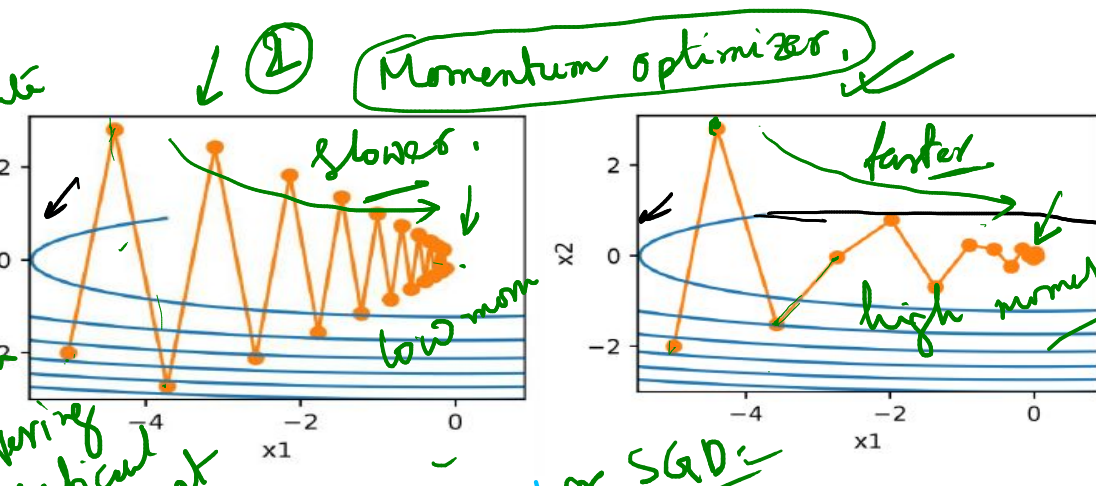
Momentum optimize.

γ low

GM

Momentum opt

High ✓

GM

a. Figure A is SGD momentum optimizer with high momentum, Figure B is RMSProp or AdaGrad and Figure C is SGD momentum optimizer with low Momentum.

b. Figure A is RMSProp or AdaGrad, Figure B is SGD Momentum with low Momentum and Figure C is SGD momentum with high momentum.
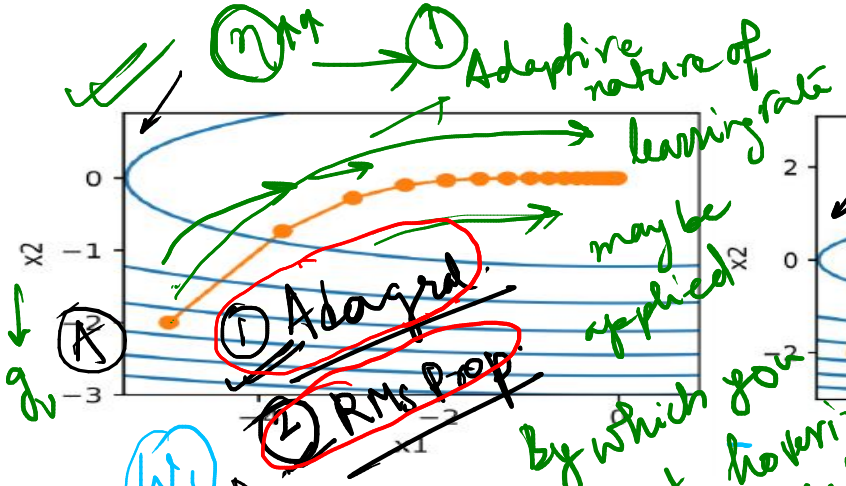
c. Figure A is SGD Momentum optimizer with low momentum, Figure B is RMSProp or AdaGrad and Figure C is SGD Momentum optimizer with high Momentum.

d. None of the above

Contour plot
Equal energy line

$\eta \uparrow$ → ① Adaptive nature of learning rate

② Momentum optimizer

may be applied

By which you are not traveling into vertical direction unlike the

slower

low mom.

faster

high moment

Ⓐ

g

① Adagrad

② RMSProp

SGD Ⓦ₁

momentum SGD Ⓦ₂

Momentum.

→ Momentum opt or SGD:
$$W_{n+1} = W_n - \eta \nabla_W L(i) + \gamma V_{n-1}$$

changes $\eta$ based on the gradient

scales the gradient

gᵥ

Ⓦ₂

$g_v > g_h$

MSGD → Vanilla

Basic version of any algorithm

$S_5$ Constant
$S_4 → S_3, S_4$
$S_1 → S_1, S_2, S_3$

Exponential gradient Avg decay

$S_t = \beta S_t + (1-\beta) S_{t-1}$

$\beta = 0.9$

$g_t, r_t,$

$\eta \uparrow \frac{g_t}{\sqrt{\epsilon + r_t}}$

$$W_{n+1} = W_n - \frac{\eta}{\sqrt{\epsilon + r_t}}$$

$\eta, \downarrow$

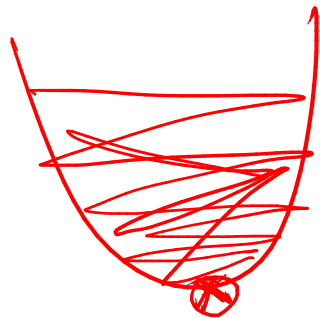What can be a possible consequence of choosing a very small learning rate? Choose the correct option.

a. Slow convergence

b. Overshooting minima $\to \eta \uparrow\uparrow$.

c. Oscillations around the minima $\to \eta \uparrow\uparrow$.

d. All of the above

$\eta \downarrow\downarrow$

① High settling time
② Slower convergence.

Two version of SGD are implemented as follows:

SGD1: SGD1 samples data points in same order for every epoch while constructing minibatch

SGD2: SGD2 samples data samples in random order for every epoch to construct minibatch

→ stochastic → Random in nature

Select the correct statement

    a. SGD1 is faster than SGD2 and robust to local minima entrapment

    b. SGD2 is faster than SGD1 and robust to local minima entrapment

    c. SGD1 and SGD2 have same convergence characteristics

    d. None of above

model. Compile(   +

Shuffle = True

Memorize the features.
X Generalize the understanding.

RMSProp resolves the limitation of which optimizer?
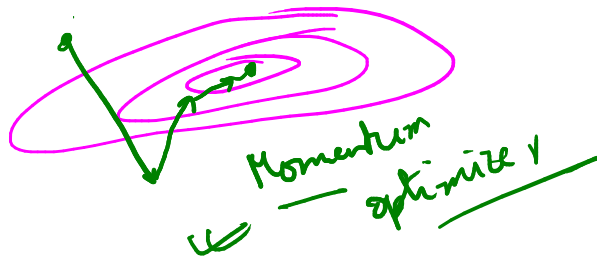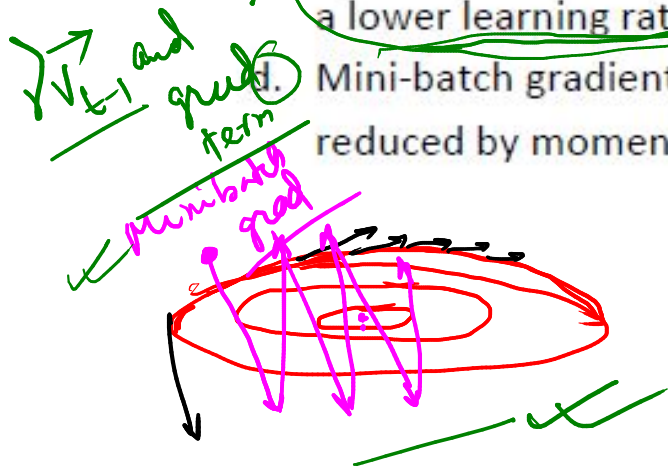
a. Adagrad (exponential decaying Avg. grad)

b. Momentum

c. Solves problem of option b but not a

d. Neither a nor b

$$r_t = \beta S_t + (1-\beta) S_{t-1}$$

Which of the following is a possible advantage of momentum optimizer over of mini-batch gradient descent?

*Vanilla*

a. Mini-batch gradient descent performs better than momentum optimizer when the surface of the loss function has a much more elongated curvature along X-axis than along Y-axis

b. Mini-batch gradient descent always performs better than momentum optimizer

c. Mini-batch gradient descent will always overshoot the optimum point even with a lower learning rate value → *only for high η (overshoot will be there)*

d. Mini-batch gradient might oscillate in its path towards convergence which can reduced by momentum optimizer

*$\vec{V}_{t-1}$ and grad term*

*minibatch grad*

*Momentum optimizer*

The following is the equation of update vector for momentum optimizer. Which of the following is true for $\gamma$?

Momentum factor / term

$$V_t = \gamma V_{t-1} + \eta \nabla_\theta J(\theta)$$

gradient associated with learning rate $\eta$

a. $\gamma$ is the momentum term which indicates how much acceleration you want
b. $\gamma$ is the step size
c. $\gamma$ is the first order moment
d. $\gamma$ is the second order moment

Velocity by which you have landed to $w_t$ from $w_{t-1}$

Why it is at all required to choose different learning rates for different weights?

a. To avoid the problem of diminishing learning rate (To bancalisen
b. To avoid overshooting the optimum point
c. To reduce vertical oscillations while navigating the optimum poin
d. This would aid to reach the optimum point faster

It ensures the convergen of the problem optimal in a faster manner.

$W_{n+1} = W_n - \eta \nabla L$

stopping of the gradient value

Let $J(\theta)$ be the cost function. Let the gradient descent update rule for $\theta_i$ be,

$$\theta_{i+1} = \theta_i + \nabla\theta_i$$

What is the correct expression of $\nabla\theta_i$, $\alpha$ is the learning rate.

$-\alpha \dfrac{dJ(\theta_i)}{d\theta_i}$
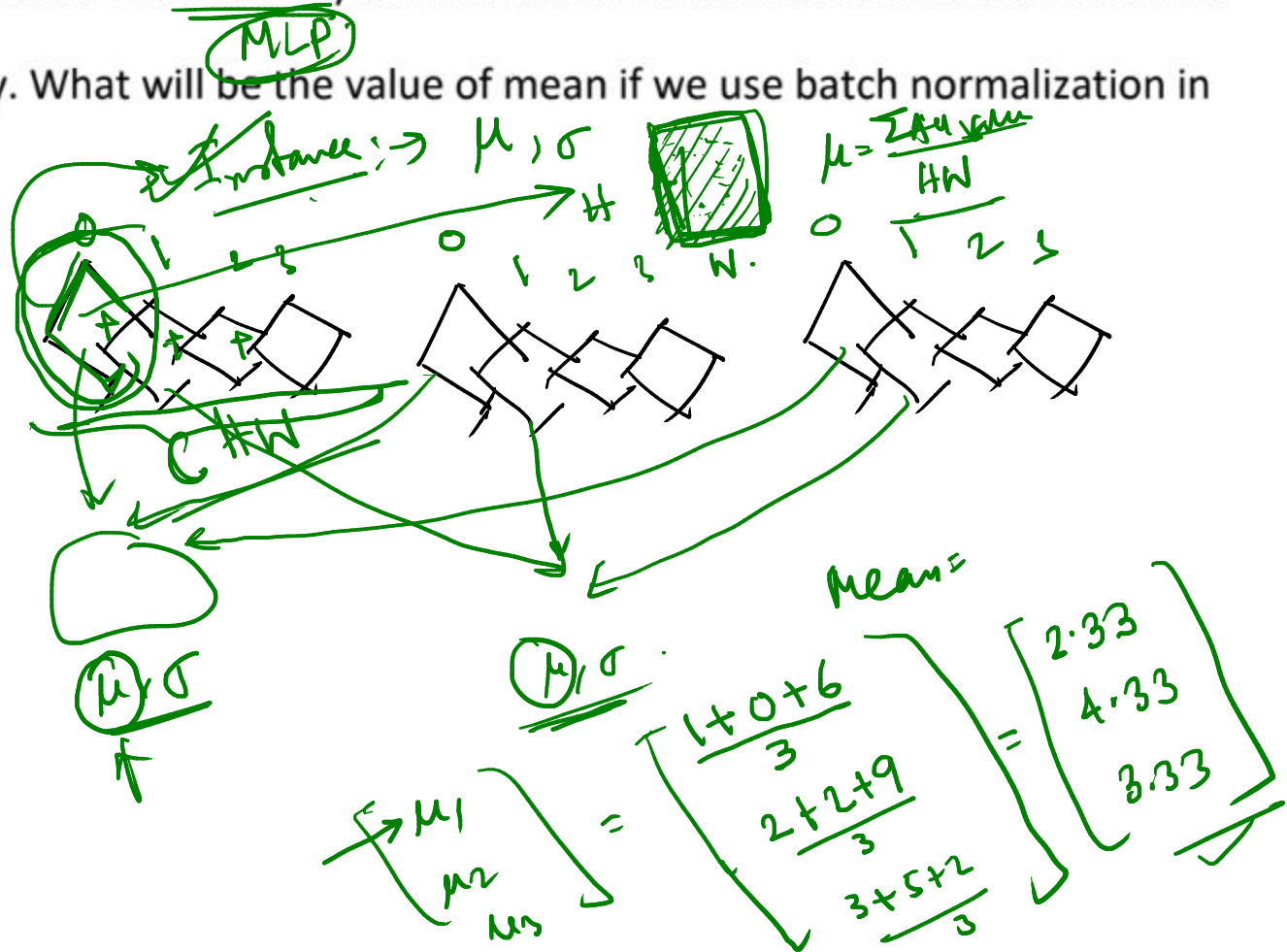
$\alpha \dfrac{dJ(\theta_i)}{d\theta_i}$

$-\dfrac{dJ(\theta_i)}{d\theta_{i+1}}$

$\dfrac{dJ(\theta_i)}{d\theta_i}$

$\theta_{i+1} = \theta_i - \alpha \cdot \dfrac{\partial J(\theta_i)}{\partial \theta_i}$
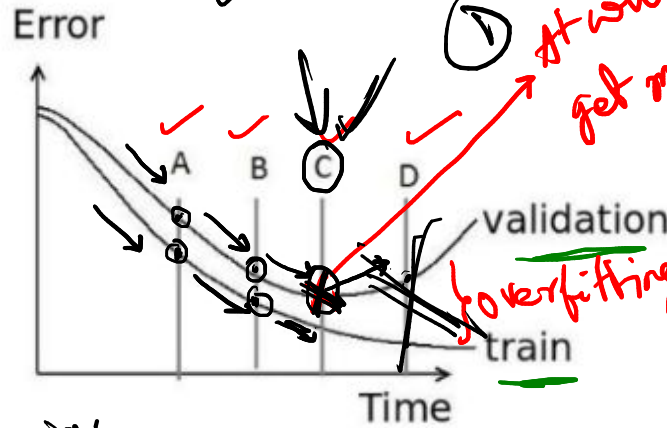
A neural network has 3 neurons in a hidden layer. Activations of the neurons for three batches are $\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 2 \\ 5 \end{bmatrix}$, $\begin{bmatrix} 6 \\ 9 \\ 2 \end{bmatrix}$ respectively. What will be the value of mean if we use batch normalization in this layer?

a. $\begin{bmatrix} 2.33 \\ 4.33 \\ 3.33 \end{bmatrix}$

b. $\begin{bmatrix} 2.00 \\ 2.33 \\ 5.66 \end{bmatrix}$

c. $\begin{bmatrix} 1.00 \\ 1.00 \\ 1.00 \end{bmatrix}$

d. $\begin{bmatrix} 0.00 \\ 0.00 \\ 0.00 \end{bmatrix}$

*(Handwritten annotations in green:)*

MLP

$\mu, \sigma$

$\mu = \dfrac{\Sigma \text{ all value}}{HW}$

Mean $= \begin{bmatrix} \dfrac{1+0+6}{3} \\ \dfrac{2+2+9}{3} \\ \dfrac{3+5+2}{3} \end{bmatrix} = \begin{bmatrix} 2.33 \\ 4.33 \\ 3.33 \end{bmatrix}$
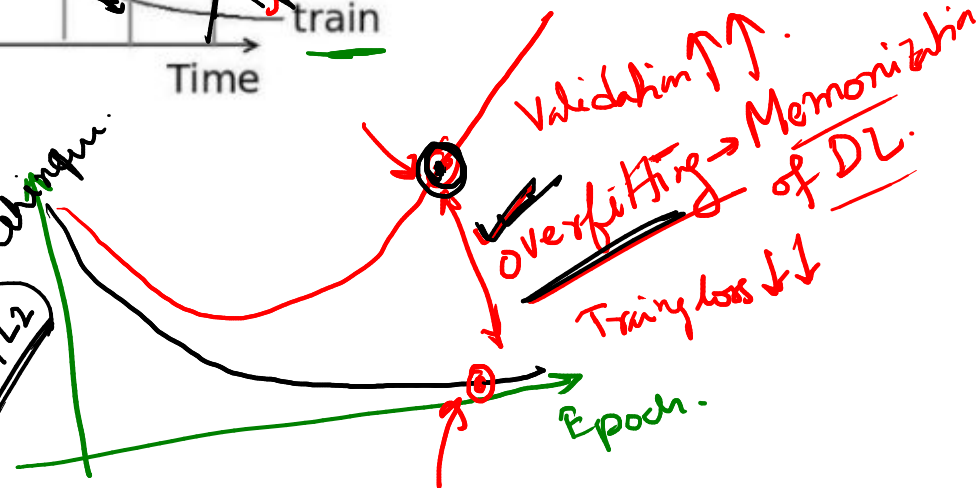
$\mu_1, \mu_2, \mu_3$

$\mu, \sigma$

C HW

2. While training a neural network for image recognition task, we plot the graph of training error and validation error. Which is the best for early stopping?



Early stopping
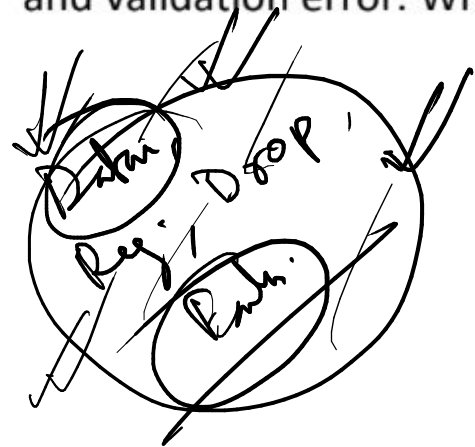↳ Premature stopping of updation of weights to alliviate the issue of overfitting.
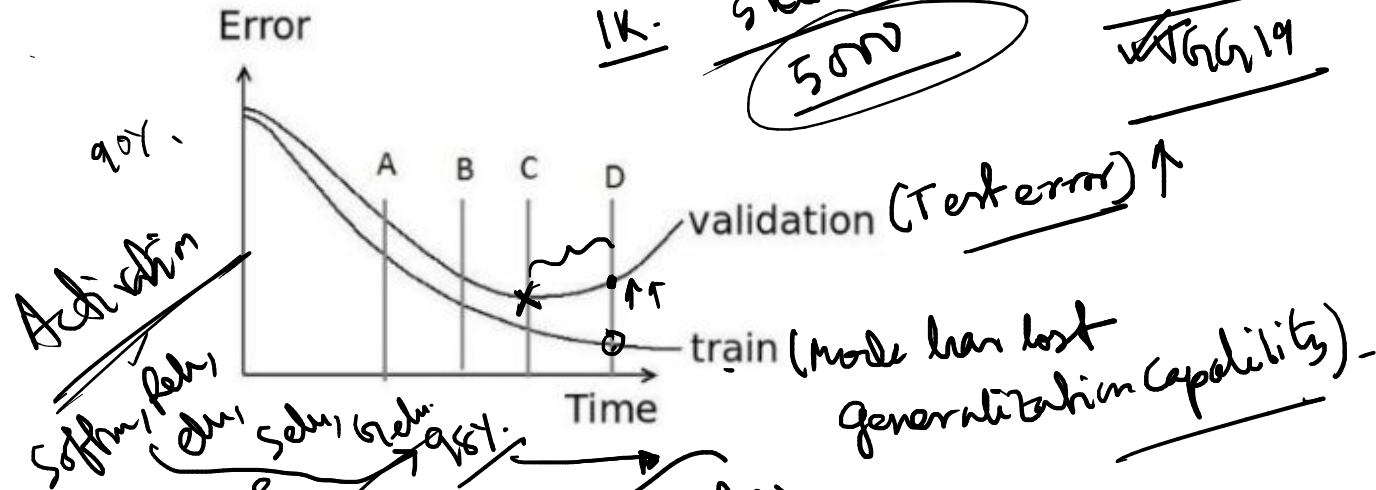
Error

A   B   C   D

validation

train

Time

① At which you get minimum validation error

② Beyon this you get overfitting

Overfitting has started

a. A
b. B
c. C
d. D

Overfitting
↳ Early stopping
↳ Dropout
↳ Regularation (d₁, L₂)

Error.

Training Techniqu.

Validation ↑↑.
Overfitting → Memorization of DL.
Traing loss ↓↓

Epoch.

While training a neural network for image recognition task, we plot the graph of training error and validation error. Which is the best for early stopping?



a. A
b. B
c. C
d. D

Error

A  B  C  D

validation (Test error) ↑

train (Mode han lost generalization capability)

Time

90%.

Activation
Softmy Relu
elu
selu
Relu
98%.

IK.  5 class
5000

√VGG 16
√VGG 19

Imagenet challenge
Alex, Vgg, Resnet, Densenet
Dense
Data loss
+ Regularization loss

CVPR, ICLR.
ConvNet

5%. NLP
Transformer
chatGPT
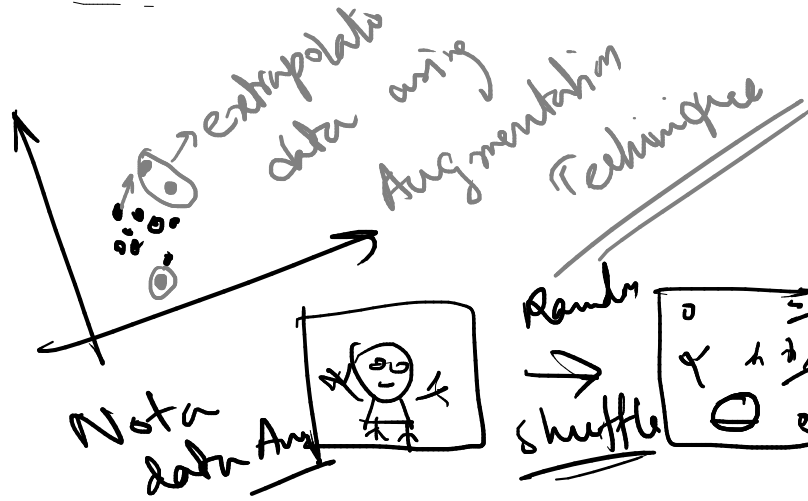
Data
Reg, Drop
Batch

Which among the following is NOT a data augmentation technique?

*Image*     d

☒ Random horizontal and vertical flip of image   Augmentation

b. Random shuffle all the pixels of an image

☒ Random color jittering ←

d. All the above are data augmentation techniques

Rotation is an augmentation Technique

→ extrapolate data using Augmentation Techniques

Not a data Aug

Random shuffle

$n_B$
$n_G$
$n_R$

R G B.
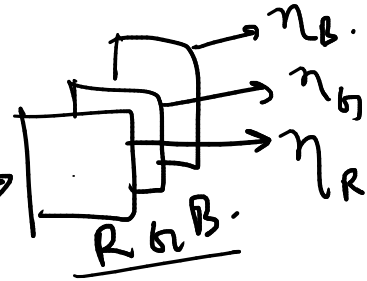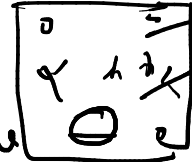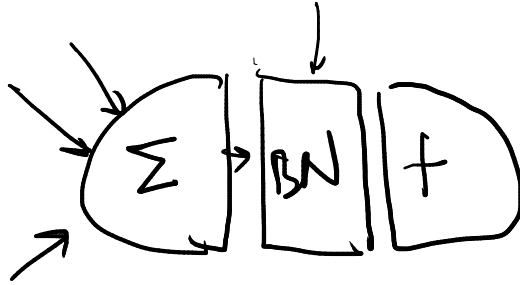
Image + Noise

Batch Normalization is helpful because

    a.  It normalizes all the input before sending it to the next layer

    b.  It returns back the normalized mean and standard deviation of weights

    c.  It is a very efficient back-propagation technique

    d.  None of these

A Batch Norm layer accepts batch of (128D) vector. How many parameters of Batch norm get trained via backpropagation during the course of training

a. 256
b. 512
c. 128
d. 1024

BN $\downarrow$ $\downarrow$
$\gamma$, $\beta$.

BN

$1 \rightarrow 2$

$128D \rightarrow 2 \times 128$

$= 256$

Which of the following is a regularization method?

    a.  Data augmentation
    b.  Dropout
    c.  Weight decay
    d.  All of the above

Two variant training schedulesamples its minibatches in the following manner
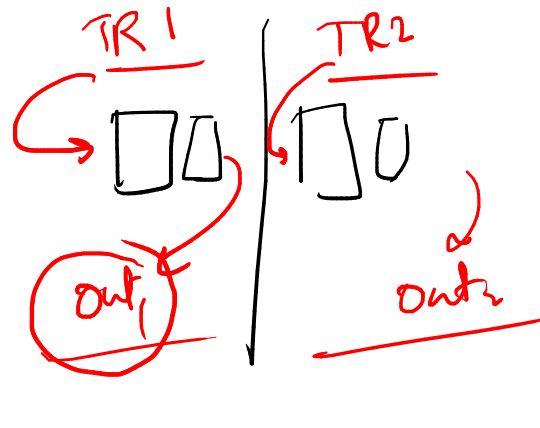
Training Schedule 1

Mini batch 1=[Image1, Image2, Image3]
Mini batch 2=[Image4, Image5, Image6]

Training Schedule 2

Mini batch 1=[Image1, Image4 , Image3]
Mini batch 2=[Image2, Image 5, Image6]

The output activations of each corresponding image is compared across Training schedule 1 and Training schedule 2 for a CNN with batch norm layers. Choose the correct statement

a. Activation outputs of corresponding image will be same across Training schedule 1 and Training schedule 2

b. Activation outputs of corresponding image will be different across Training schedule 1 and Training schedule 2

c. Some activations outputs of corresponding images will be same but some will be different

d. None of these.

# Thank You