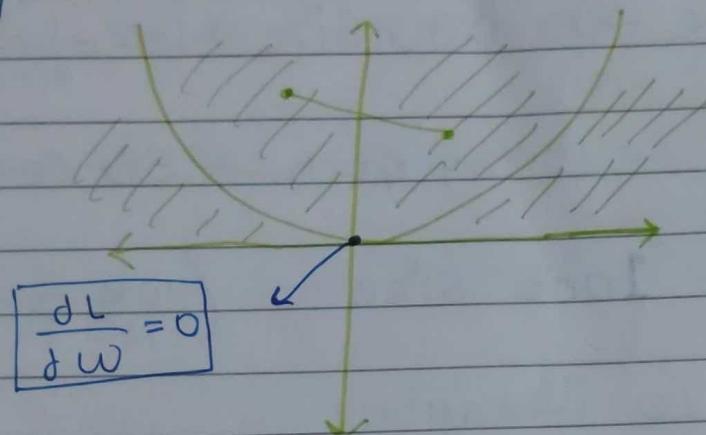
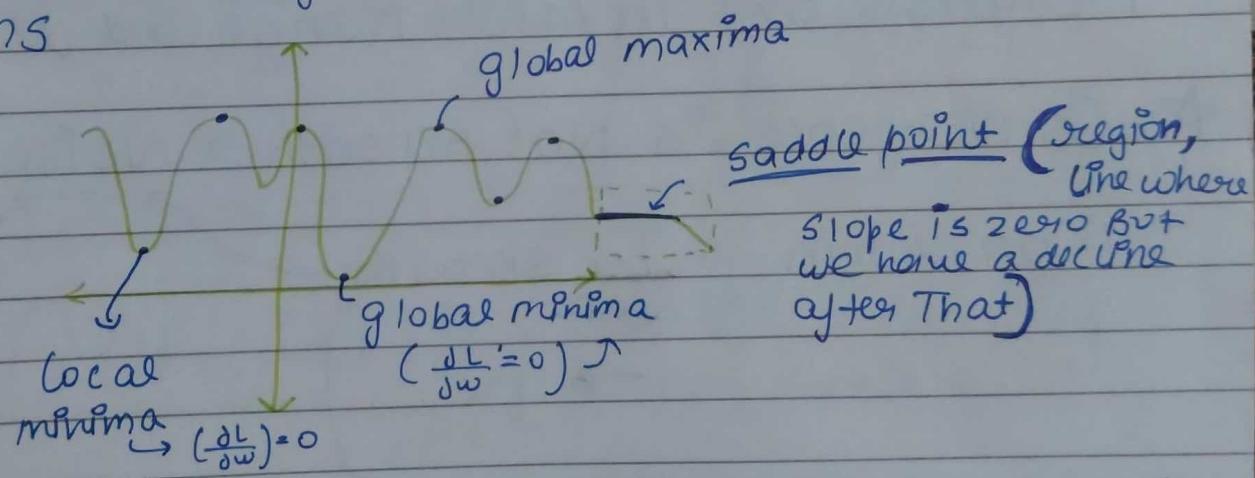


In ML → most/all the times we get Convex graph



[2 points
conn, lies
in same
Region]

In DL → most of the time we get Non-convex graphs



Because of Diff loss func and so many neurons, when weight has to be optimised for every point.

* Derivative $\rightarrow \frac{\partial \text{Loss}}{\partial w_{\text{old}}}$ at minima point is zero

\hookrightarrow it means $w_{\text{new}} = w_{\text{old}}$ and we now know we have reached an optimal soln.

Problem 2

In DL-Curve \rightarrow the derivative at local minima is zero but we haven't reach the global minima

Solution \rightarrow Advance optimisers (apart from GD, SGD, mSGD) and Adv Loss func.

→ SGD with Momentum

$(K < n)$

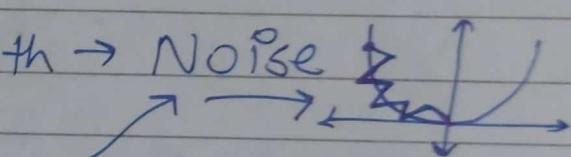
- In MB SGD → we take Batch Size → loss = $\frac{1}{K} \sum_{i=1}^K (y - \hat{y})^2$

↳ 10K records so every epoch we set same size
 Epoch = 1000 size $(\frac{10K}{1K} \Rightarrow 10)$

1 Epoch = 10 iteration.

20 Epoch = 200 iterations

- Iterations are in mid-range and lil faster.

★ Convergence is Not smooth → Noise 

↳ Momentum → To smooth

$$\omega_t = \omega_{t-1} - \eta \frac{dL}{d\omega_{t-1}}$$

We will calculate → Exponential Weighted Average

- EWA → Calculating the 'new value' on the basis of importance of previous value.

→ $t_1, t_2, t_3, t_4, \dots, t_n$ Iteration
 → $a_1, a_2, a_3, a_4, \dots, a_n$ value

$$V_{t_1} = a_1$$

$$V_{t_2} = \beta \cdot V_{t_1} + (1-\beta) a_2$$

$$V_{t_3} = \beta \cdot V_{t_2} + (1-\beta) a_3$$

β - importance param
 if $\beta = 0.9$ ($a_1 \uparrow$ and $a_2 \downarrow$)
 if $\beta = 0.1$ ($a_1 \downarrow$ and $a_2 \uparrow$)

EWA

$$\rightarrow w_t = w_{t-1} - \eta v_{dw}$$

$$\hookrightarrow v_{dw_t} = \beta v_{dw_{t-1}} + (1-\beta) \frac{\delta L}{\delta w_t}$$

$\Rightarrow \beta = 0.95$

$\text{Input} \rightarrow v_{dw} = 0.$

$$v_{dw_{t-1}} = \beta v_{dw_{t-2}} + (1-\beta) \frac{\delta L}{\delta w_{t-1}}$$

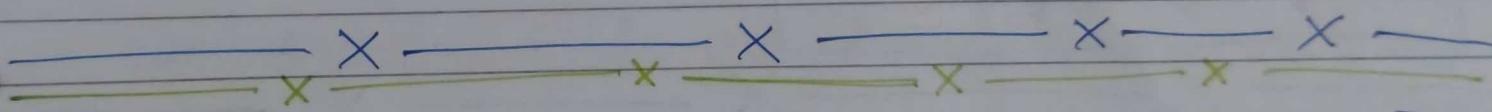
$$\frac{\delta L}{\delta w_{old}}$$

~~Bias~~

$$w_{new} = w_{old} - \alpha \frac{\delta L}{\delta w_{old}}$$

$$B_{new} = B_{old} - \alpha \frac{\delta L}{\delta B_{old}}$$

Similar steps for weights and bias



Yth] ADAGRAD [Adaptive Gradient Boosting]

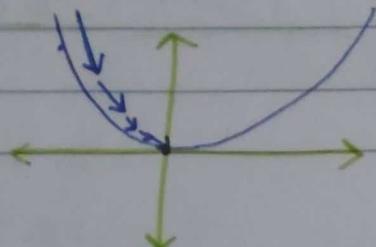
- Because of this, we stopped using EWA

- 1) here our learning Rate is NOT FIXED

- 2) our α keeps changing itself as the time goes in training

- 3) at first \rightarrow Big Jump and then small ones
 $\alpha \rightarrow$ change dynamically

$$\alpha = 0.5 \rightarrow 0.4 \rightarrow 0.3 \rightarrow 0.1 \rightarrow 0.01 \leftarrow 0.05 \leftarrow 0.09$$



- $t = \text{new} \rightarrow \underline{\text{curr iter}}$
- $t-1 = \text{old} \rightarrow \underline{\text{prev iter}}$

$$\text{GD} \rightarrow W_t = W_{t-1} - \alpha \frac{\delta L}{\delta W_{t-1}}$$

$$\text{AGD} \rightarrow W_t = W_{t-1} - \alpha' \frac{\delta L}{\delta W_{t-1}}$$

$$\alpha' = \frac{\alpha}{\sqrt{n_t + \epsilon}} \quad \begin{matrix} \text{initial value} \\ \left[\begin{matrix} \text{Small value number} \\ \text{if it is zero} \end{matrix} \right] \end{matrix}$$

If at $t_1 = 0.001 \rightarrow t_2 = 0.009 \rightarrow t_3 = 0.001, \rightarrow t_4 = 0.005$

To $\downarrow \alpha'$, $n_t \uparrow (\epsilon)$

PROBLEM \rightarrow In deep Multi-H Neural, after some time the value of n_t will be Huge and α' will be very very small \rightarrow Vanishing Prob

due to Learning Rate \leftarrow NO Convergence

SOLUTION \rightarrow

5th] Ada Delta and Rms Prop

Recuse

$$w_t + \bar{v}$$

$$\begin{array}{l} A^L + \text{NLP} \\ A^L + \text{pro ratio} \end{array}$$

loss function

$$f = m \otimes b$$

$$0.46 \rightarrow 0 \xrightarrow{\text{reg}} 0.54$$

(restricting large values)

$$\rightarrow \alpha' = \frac{\alpha}{\sqrt{S_{dw} + \epsilon}}$$

{ Initially $S_{dw} = 0$
Just like EWA / SGD moment }

$$\text{Moving Average} = S_{dw_t} = \beta S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w} \right)^2$$

$$S_{dw} = \beta S_{dw_{t-1}} + (1-\beta) \left(\frac{\partial L}{\partial w_t} \right)^2$$

* Smaller value to restrict larger value of squared values.

$$w_t = w_{t-1} - \alpha' \frac{\partial L}{\partial w_{t-1}}$$

$$\beta_t = \beta_{t-1} - \alpha' \frac{\partial L}{\partial \beta_{t-1}}$$

* On iteration t in mini Batch

Compute $\frac{\partial L}{\partial w}$, $\frac{\partial L}{\partial \beta}$ on current mini Batch

$$S_{dw} = \beta S_{dw_{t-1}} + (1-\beta) \frac{\partial L}{\partial w_t}$$

$\beta \stackrel{\text{def}}{=} 0.95$

$$S_{dw \text{ init}} = 0$$

$$S_{db} = \underbrace{\beta S_{db_{t-1}}}_{\text{init}} + (1-\beta) \frac{\partial L}{\partial b_t}$$

$$w_t = w_{t-1} - \alpha' \frac{\partial L}{\partial w_{t-1}}$$

$$b_t = b_{t-1} - \alpha' \frac{\partial L}{\partial b_{t-1}}$$

everytime take mini-Batch of data

* 6] ADAM OPTIMIZER
{ Adaptive Moment Estimation }

Corrected

estimating η_t value with help of
 $S_{dw} \rightarrow \underline{C_{WA}}$

(smooth)

\rightarrow RMS prop + Momentum | combine - V_{dw} , V_{db}
 S_{dw} , S_{db}

1] $V_{dw} = \boxed{V_{db}} V_{db} = S_{dw} = S_{db} = 0$

on iteration t]

2] Compute $\frac{dL}{dw}$, $\frac{dL}{db}$ using current mini Batch

3] $\boxed{V_{dw,t} = \beta_1 V_{dw,t-1} + (1-\beta_1) \frac{dL}{dw,t}}$

$\boxed{V_{db,t} = \beta_1 V_{db,t-1} + (1-\beta_1) \left(\frac{dL}{db,t} \right)}$

$\boxed{S_{dw,t} = \beta_2 S_{dw,t-1} + (1-\beta_2) \left(\frac{dL}{dw,t} \right)^2}$

$\boxed{S_{db,t} = \beta_2 S_{db,t-1} + (1-\beta_2) \left(\frac{dL}{db,t} \right)^2}$

4] $\boxed{w_{old/t} = w_{t-1} - \alpha' V_{dw}}$

or ↴

$w_t = w_{t-1} - \eta' V_{dw}$

$w_t = w_{t-1} - \frac{\alpha \cdot V_{dw}}{\sqrt{S_{dw,t} + \epsilon}}$

$w_t = w_{t-1} - \frac{\eta \cdot V_{dw,t-1}}{\sqrt{S_{dw,t-1} + \epsilon}}$

$b_t = b_{t-1} - \frac{\eta \cdot V_{db,t-1}}{\sqrt{S_{db,t-1} + \epsilon}}$

• Bias Correction •

$\rightarrow V_{dw}^{\text{corrected}} = \frac{V_{dw}}{(1-\beta_1^t)}$

$V_{db}^{\text{corrected}} = \frac{V_{db}}{(1-\beta_1^t)}$

$$S_{dw}^{\text{correction}} = \frac{S_{dw}}{(1-\beta_2^t)} + S_{db}^{\text{corrected}} = \frac{S_{dw}}{(1-\beta_2^t)}$$

New formula

- $w_{\text{new}} = w_{\text{old}} - \frac{\alpha}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon} \cdot \sqrt{v_{dw}}$

$$w_t = w_{t-1} - \frac{\alpha}{\sqrt{S_{dw_{t-1}}^{\text{corrected}}} + \epsilon} \cdot \sqrt{v_{dw_{t-1}}}$$

- $\beta_t = \beta_{t-1} - \frac{\alpha}{\sqrt{S_{db_{t-1}}^{\text{corrected}}} + \epsilon} \cdot \sqrt{v_{db_{t-1}}}$

DIFF. TYPES OF LOSS FUNCTION

1) Regression \rightarrow Squared error loss (single point)

$$L_d = (y - \hat{y})^2$$

• Mean squared error \rightarrow [MSE] $\frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 = T$

• Absolute Loss error \rightarrow [MAE] $\frac{1}{m} \sum_{i=1}^m |y_i - \hat{y}_i|$

★ MAE is more Robust to outliers as compared to MSE.

(3) Huber loss $\rightarrow \text{MSE} + \text{MAE}$

$$\rightarrow \begin{cases} \frac{1}{2}(y - \hat{y})^2, & \text{if } |y - \hat{y}| \leq \delta \\ \delta|y - \hat{y}| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$$

Hyperparameter

- \rightarrow Fix all prob of MSE and MAE.
- \rightarrow Prob of outliers is fixed

— X — X — X — X —

• CLASSIFICATION

\rightarrow ① cross Entropy (in log Reg usecase)

$$\mathcal{L} \Rightarrow -y * \log(\hat{y}) - (1-y) * \log(1-\hat{y})$$

\Downarrow

$y = 0/1$

Binary

cross

Entropy

$$\begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log(\hat{y}) & \text{if } y=1 \end{cases}$$

$$\hat{y} = \text{Sigmoid.} \rightarrow \frac{1}{1+e^{-z}} \quad \mid \quad \frac{1}{1+e^{-y}}$$

Categorical

O/p → Vector → one Hot Encoding

• MULTI-CLASS CROSS ENTROPY LOSS

$$L(x_i, y_i) = - \sum_{j=1}^c y_{ij} \cdot \log(\hat{y}_{ij})$$

→ Perform
one-Hot -①
Encoding.
of output feat.

• O/p → multi-class → O/p layer

↳ Good, Bad, Neutral → One hot encoding

↳ Good : Bad, Neutral
[0 : 0 : 1]

$f_1 \ f_2 \ f_3 \ \dots \ O/p$

↓ one hot [O/p → unique features → own column]

$$\begin{matrix} f_1 & f_2 & f_3 & y_1 & y_2 & y_3 & [y_{ij}] \\ \dots & \dots & \dots & [1, 0, 0] \\ \dots & \dots & \dots & [0, 1, 0] \\ & & & y_{21} & y_{22} & y_{23} \end{matrix}$$

$i = \text{Row number}$
 $j = \text{No. of Bits} \rightarrow \text{Columns}$
 $c = \text{Categories/Bits}$

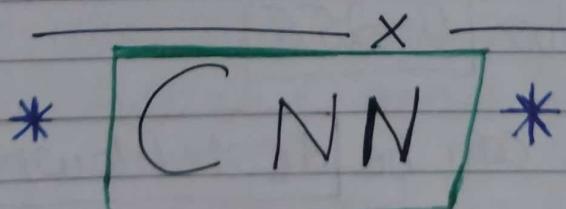
• We use SOFTMAX act. func for multi-class problems.

$$\sigma(z) = \frac{e^{z_i}}{\sum_{j=1}^c e^{z_j}} \rightarrow [P_1, P_2, P_3, \dots, P_c] \rightarrow \text{majority wins.}$$

↳ Prob value

• To calculate $\hat{y}_{ij} \rightarrow \text{SOFTMAX} \rightarrow \frac{e^{10}}{e^{10} + e^{20} + e^{30} + e^{40}} \rightarrow \text{Probability}$

↳ Backtrack ← calculate ← Put in Loss ←

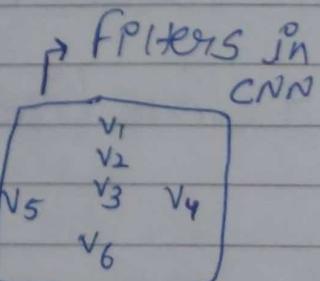


HUMAN BRAIN

image → eyes → Brain

Cerebral Cortex

↳ Visual Cortex

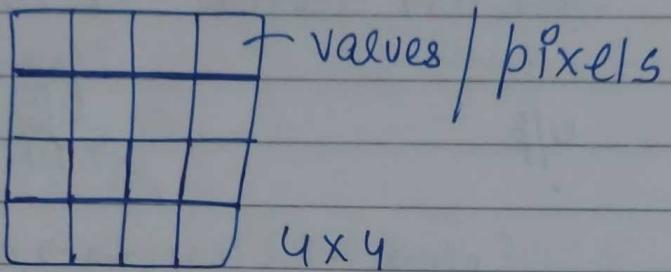


↳ Various layers that detect diff. feat of images

→ Convolution in CNN.

→ Basics About Gray scale Image

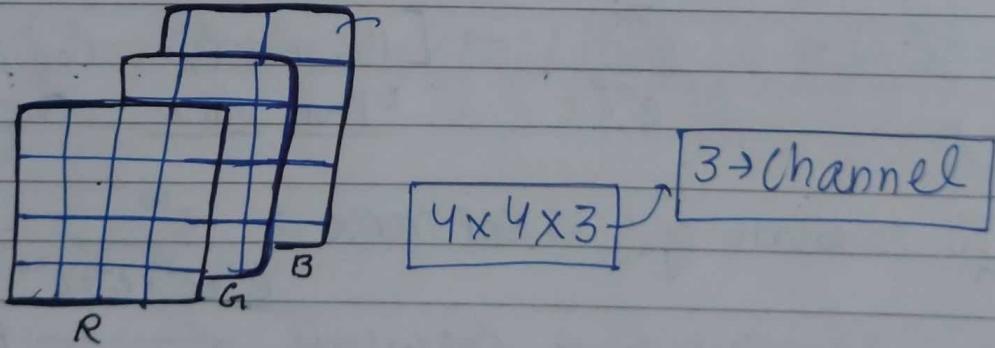
- Define as $a \times a$ [eg. $\rightarrow 4 \times 4, 6 \times 6$ etc..] [pixels \times pixels]



Values Range B/w $0 - 255$

→ Basics About RGB | Color Image

- Stacked up windows / matrix / matrices of (values/pixels)

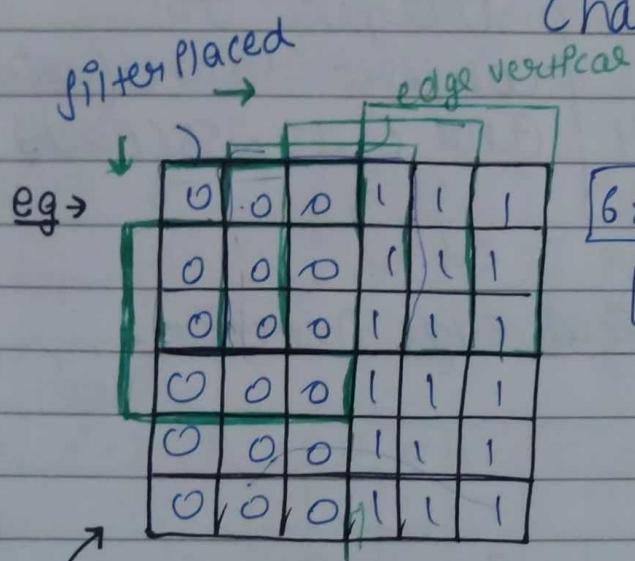


Values Range B/w $0 - 255$

* These pixel values can be scaled down.

Normalization \longleftrightarrow Standardization

* Convolution: Putting filter over the channel. and moving it by 1 (\rightarrow and \downarrow)



$6 \times 6 \rightarrow$ Gray scale

↳ 0 - 255

↳ Normalized

↳ 0 (white)

↳ 1 (black)

Filter → layer

↳ To detect edges / things in a image.

Here we use Filter

filter

↳ Vertical edge Detector → To detect vertical edge

1	0	-1
2	0	-2
1	0	-1

filter → vertical edge detector

• Generalized STEPS for ALL FILTERS

So STEP 1 → Place filter over the image pixel window

So STEP 2 → multiply respective cell values and apply Σ on it / them

So STEP 3 → Use that value as a cell of the result channel / window

So STEP 4 → Move / Stride one step left and Repeat STEPS 1 to 3. Once hit Boundary move 1 step down and Repeat STEPS 1 to 3. [STRIDE AND PADDING]

FINAL STEP

Output \rightarrow

0	-4	-4	0
0	-4	-4	0
0	-4	-4	0
0	-4	-4	0

4x4

Normalized the Output to
Convert -4 to 0 (white)
and 0 to 255 (black)

255	0	0	255
255	0	0	255
255	0	0	255
255	0	0	255

vertical edge detection.

* Filter for Horizontal edge Detection will look like $\rightarrow [3 \times 3] \rightarrow$

1	2	1
0	0	0
-1	-2	-1

Transpose of vertical edge detector

— X — X — X — X — X —

Formula $\rightarrow [n-f+1] \quad (6-3+1 \rightarrow 4)$

We loss some info about the image if we stride by 1 step. (from 6×6 to 4×4).

— X — X — X — X —

To prevent it from happening we use the concept of PADDING

Let's say \rightarrow we want $[6 \times 6 \rightarrow 6 \times 6]$ (No loss)

$$\text{So } n-f+1 = 6 = n = 6+3-1 = n = 8$$

So we have to convert our 6×6 to 8×8 in order to prevent loss of data.

- PADDING → converts our image to desirable dimension so to prevent any loss of data during convolution.

Padding
 $\rho = 1$

0	0	6	×	6	0	0	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	1	1	1	0
0	0	0	0	0	0	0	0

To fill
Values in
New col
and rows
↳ O Padding
↳ Neighbour Padding

- 0 Padding → Fill 0's everywhere in new cells
 - Neighbour Padding → Fill the nearest value in the new cell (comp to values in original image).

$$n-f+1 \rightarrow \text{After Padding} \rightarrow n+2P-f+1$$

- Inf loss \rightarrow $n-f+1$ [Stride] less edges
- No Inf loss \rightarrow $n+2p-f+1$ More edges.
[Stride + Padding]

Backpropagation $\xrightarrow{\text{in CNN}}$ update weights/value of filters.

• Operation of CNN (Backpropagation and filter update)

• filter = kernel

• Multiple filters

• Just like ANN $\rightarrow z = w_1 x_1 \rightarrow z = \sigma(z)$

in CNN \rightarrow Image $\square \star$ filter $\square \rightarrow$ output

$n - f + 1$

Output = $\sigma(z)$

new values
of output

→ Same (loss func \rightarrow optimiser \rightarrow backtracking \rightarrow weight upd)

happens in CNN also.

Value
updation

X — X — X — X —

MAX POOLING

→ either we can stack multiple convolution operations or stack multiple filters, or we can also apply MAX POOLING.

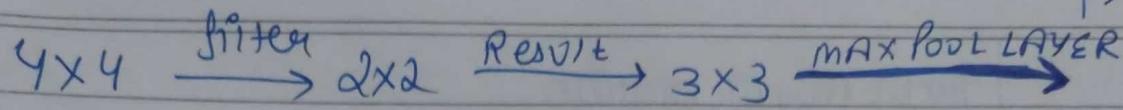
* Location Invariant \rightarrow some neurons get automatically updated when they see an Image.

* MAX POOLING \rightarrow A layer or kernel or Tensor (n-dimensional matrix/array) applied after 1/multiple convolutional operations.

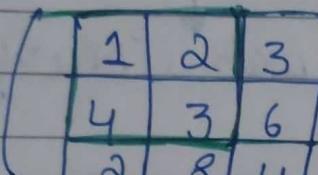
$$\begin{aligned} n - f + 1 \\ 4 - 2 + 1 \Rightarrow 3 \end{aligned}$$

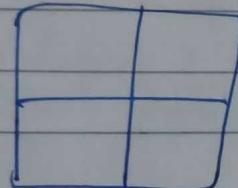
$(P=0) \rightarrow \text{Inj loss}$

Let's $\Rightarrow 2 \times 2$



- In max pool \rightarrow stride = 2 (or [Stride = size of maxpool layer])

max  \rightarrow Apply max pool layer and consider the maximum value from that Tensor or Kernel.



\rightarrow maxpool of 2×2 , so stride = 2

Output \rightarrow

4	6
8	4

max intensity value
is picked from the
convolution operation.
 \hookrightarrow So proper detection.

- These can be stacked. (Horizontally or Vertically)
- Some ex \rightarrow min pooling / avg pooling / max pooling.

DATA AUGMENTATION

- Transforming the same Image / Data into Diff forms.
- flip \rightarrow Horizontal shift \rightarrow Vertical shift \rightarrow zooming \rightarrow ADD NOISE.
- Why? \rightarrow To \uparrow variance in data so that our model performance can \uparrow and can be Robust to any way the data is presented to it.

Code snippet

Data Augmentation

• from keras.preprocessing.image import ImageDataGenerator

↓
main library.

• from tensorflow.keras.utils import img_to_array, array_to_img, load_img

• datagen = ImageDataGenerator(\hookrightarrow #parameters)

rotation_range = 40,

width_shift_range = 0.2,

height_shift_range = 0.2,

shear_range = 0.2,

zoom_range = 0.2,

horizontal_flip = True,

fill_mode = 'nearest')

(: img = load_img("Directory towards Image")

\rightarrow x = img_to_array(img) \rightarrow R, G, B \rightarrow 3D \rightarrow 3 channels
↓ make it 4D

\rightarrow x = x.reshape((1,) + x.shape)
(1, 3D size).

flow generates batches of randomly Transformed images and saves the result to the specific directory.

• for batch in datagen.flow(x, batch_size=1, save_to_dir = " ", save_prefix = " ", save_format = " ") ..

Transfer Learning → Pre-trained model → our model.

- Using Pre-trained Models to inc the performance of our Model.
- The Pre-trained models are trained on millions of images so it ↑ the accuracy of our Model.
- We can use a pretrained CNN, one trained on millions of images, as start of our Model. This will allow us to have a very good convolution base before adding our own dense layered classifiers.

* Fine tuning → Applying Transfer learning, we often want to tweak the final layers in our conv base to work better for our problem.

- Only adjusting the final few layers → Does not affect the prev layers
- To look for only features relevant to our very specific problem.

TL → Pre-trained Model → Our Model →
fine Tuning

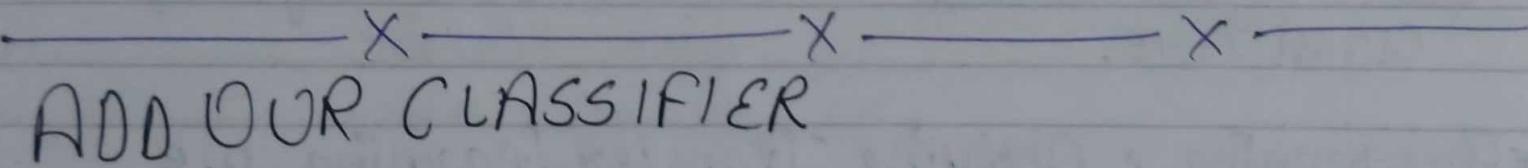
[After Training the Base Model]

Freezing The Base

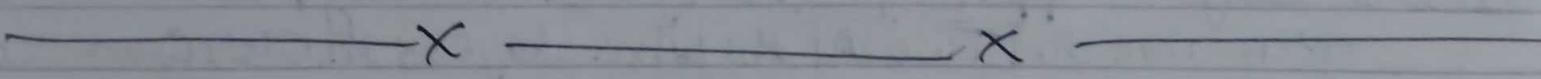
Do not train the layers of
Base model

The term freezing refers to Disabling
the training property of a layer

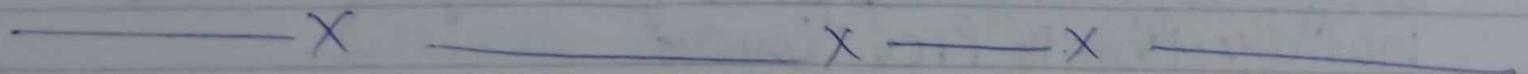
- we won't make any changes to the weights of any layers that are frozen during Training.
- Imp AS WE DON'T WANT TO CHANGE THE CONVOLUTIONAL BASE THAT ARE ALREADY HAS LEARNED WEIGHTS



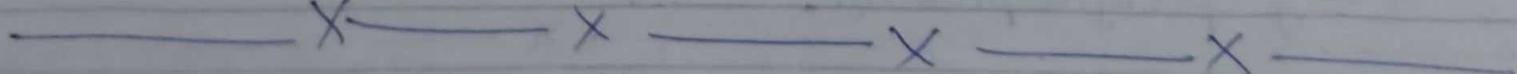
ADD OUR CLASSIFIER



Binary class → binary_crossentropy (loss)



Multiclass → Sparse-Categorical-Crossentropy (loss)



Basic ANN Prep

- Always Normalized / Standardised X-train and X-test.
- Import tensorflow as tf.

model = tf.keras.models.Sequential()

Input layer ↴

↳ optional

model.add(tf.keras.layers.Dense(units=10, kernel_initializer="he_normal", activation="relu", input_dim=13))

model.add(tf.keras.layers.Dropout(value)) → ①

Further layers - (N-layers)

↳ model.add(tf.keras.layers.Dense(units, activation))
— ① (use ①)

Output layer

↳ model.add(tf.keras.layers.Dense(units, activation))
— X — X — X — ↳ sigmoid/softmax

model.compile(optimizer, loss, metrics).

— X — X — X — ↳ optional

model.fit(X-train, Y-train, validation_split=0.30, batch_size, epochs=30)

— X — X — X —

Basic CNN Prep

Normalize img array $\rightarrow 1/255.0$

X-train and X-test \rightarrow ①

`reshape([batch], size)`

* ② eg $\rightarrow [28, 28, 1]$

model = `tf.keras.models.Sequential()`

model.add(`tf.keras.layers.Conv2D(filters, kernel_size, activation, input_shape)`)

Further layers \rightarrow model.add(`tf.keras.layers.MaxPool2D(size)`)

model.add(`tf.keras.layers.Conv2D(filters, size, activation)`)

L Similar to Dropout

Use ①

$x \longrightarrow x \longrightarrow x \longrightarrow x$

model.add(`tf.keras.layers.Flatten()`) \rightarrow makes 1D for ANN

$x \longrightarrow x \longrightarrow x \longrightarrow$

Add ANN layer \rightarrow Dense and Dropout and output layer

Same `model.compile` and `model.fit`

$x \longrightarrow x \longrightarrow x \longrightarrow x$

* Hyperparameter tuning Using Keras-tuner

! pip install Keras-tuner - Upgrade

in layers we add `hp`

CNN

`filters = hp.Int("name", min_value, max_value, step)`

`kernel_size = hp.Choice("name", values=[a, b])` (axa or bxh)

activation and input_shape remains same

$x \longrightarrow x \longrightarrow x$

ANN

`units = hp.Int("name", min_value, max_value, step)`

activation remains same

$x \longrightarrow x \longrightarrow x$

From Keras-tuner import RandomSearch

from Keras-tuner.engine.hyperparameters import HyperParameters

`ts = RandomSearch(func_name, objective='val-accuracy', max_trials=5, directory=" ", project_name)` *

`ts.get_best_models`

`ts.search(X_train, Y_train, epochs, validation_split)` (num-models=1) [o] (last step)