

DL

Epoch :- number of times data traversed for optimizing parameters.

Hidden layers minimize the error function & increases optimization & extract features from the linear classifier inputs.

Q What does a single neuron do?

* adding non trivial non linearity

If we want linear separation, there is no need of neurons or separate layers.

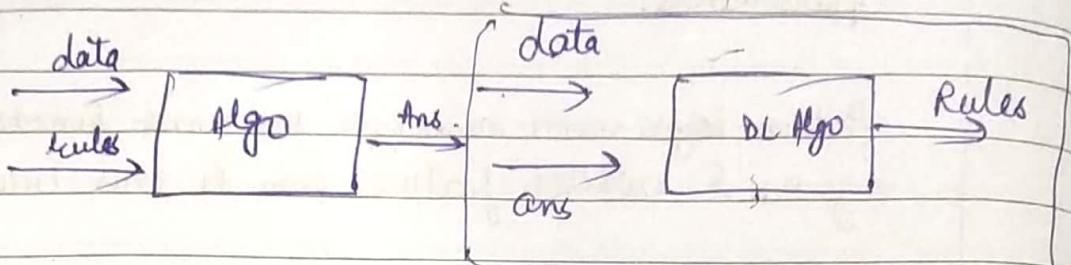
Test loss → (misclassification)

Training loss → less than test loss (misclassification).

Training set →

Image Net Data set :-

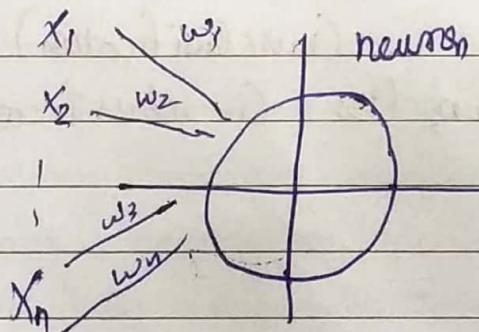
Tensorflow Specialization



Given the data & the answer, we try to find the pattern.

- * representation of data in such a manner that it learns & it learns function efficiently. (more the data better)
- * we talk mostly of supervised learning prediction of function

What does a neuron do?



$$y = wX$$

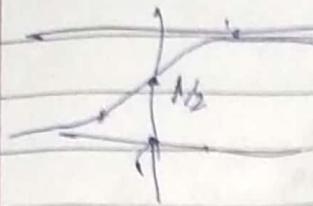
$$y = \sum_{i=1}^n x_i w_i$$

$$y = f \circ g \circ h(x)$$

$$y = \sigma \left(\sum_{i=1}^n x_i w_i \right)$$

Sigmoid f. with domain $(-\infty, \infty)$

with range $(0, 1)$



Keywords

model :-

Sequential :- put all layers sequentially (i.e. all layers one after other).

import tensorflow as tf

import numpy as np

from tensorflow import keras.

model = tf.keras.Sequential([keras.layers.Dense(units=1, input_shape=[1])])

it runs till it reaches a minimum E.

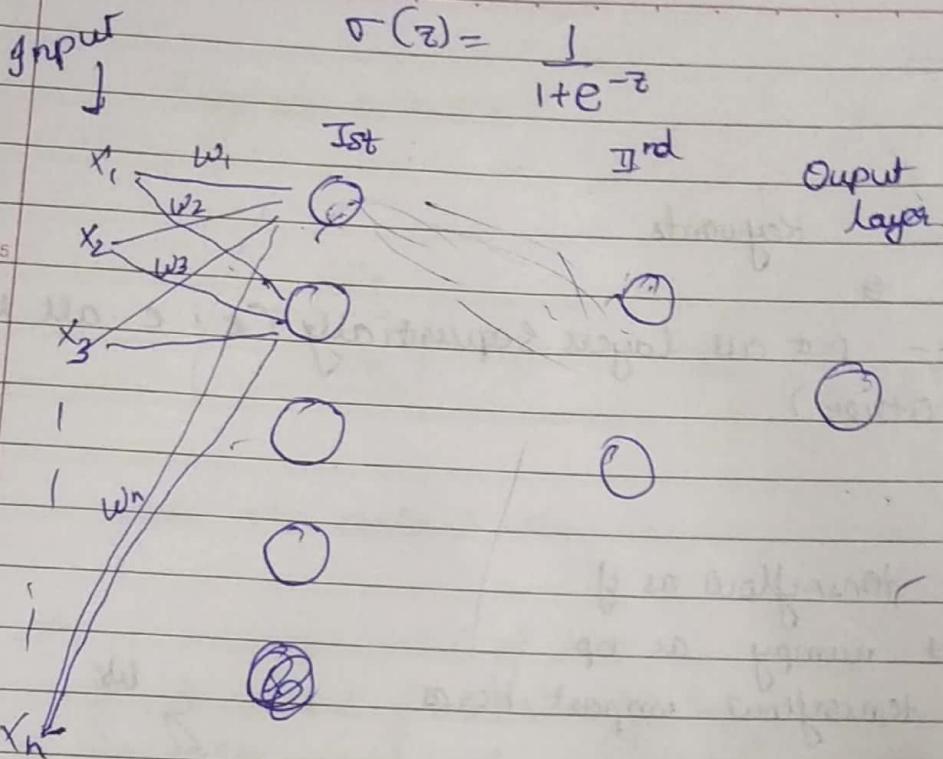
what kind of input shape receiving

model.compile(optimizer='sgd', loss='mean_squared_error')

model.fit(x, y, epochs=10)

print(model.predict([1]))

Stochastic Gradient descent is used to minimize the cost function, which is calculated on the basis of input features. Main aim is to extract important features.



If there are n neurons & m data points - then there are $4n$ parameters.

$$W^{[1]} = B$$

this represents $[0-1]$ matrix $n \times m$ where
layer. m is number of neurons

now matrix size for a general l^{th} layer would be.

$$W^{[l]} = [\quad]_{n^{[l]} \times n^{[l-1]}}$$

\downarrow \rightarrow no. of neurons
number of neurons in $(l-1)^{\text{th}}$ layer
in l^{th} layer

Every layer will have input (weights) $w^{[i]}$
 & number of neurons $n^{[i]}$, & or two outputs,

& one denoted by z & the other one by a

$$z_3^{[1]}, a_3^{[1]}$$

→ fast hard

In third layer first neuron

In third first layer
third neuron.

For a single net layer, weight matrix is represented as

To create a weight matrix for each neuron in a layer,
we will create a vector.

Output of each layer will be defined as.

(linear output i.e. z)

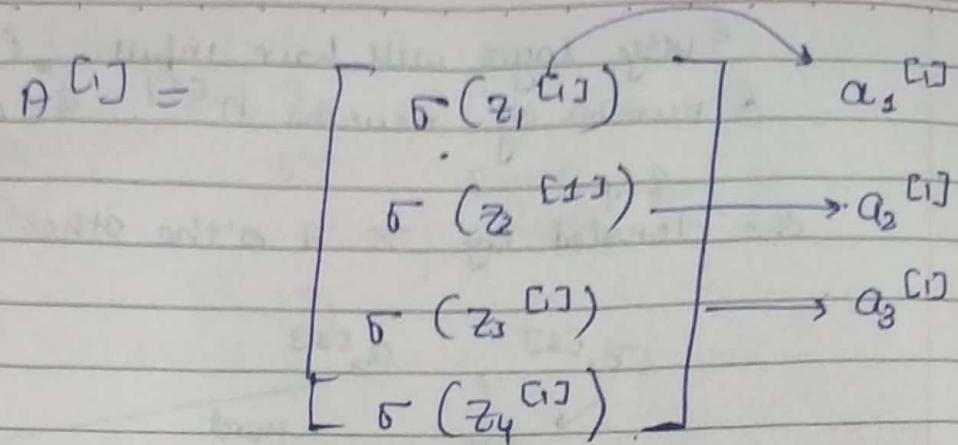
for a

single
neuron

$$\rightarrow \begin{bmatrix} w_{11}^{[1]} & \dots & w_{1n}^{[1]} \\ w_{21}^{[1]} & \dots & w_{2n}^{[1]} \\ w_{31}^{[1]} & \dots & w_{3n}^{[1]} \\ w_{41}^{[1]} & \dots & w_{4n}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

←
Output from
first layer.



~~Now~~ Now output of first layer is input of second layer

Now z , A & W for second to first to second layer.

$$w^{[2]} = \begin{bmatrix} w_{11}^{[2]} & \dots & w_{14}^{[2]} \\ w_{21}^{[2]} & \dots & w_{24}^{[2]} \end{bmatrix} \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

(linear output) $z^{[2]} = w^{[2]} \cdot A^{[1]}$

$$A^{[2]} = \sigma(z^{[2]})$$

$$z^{[l]} = w^{[l]} A^{[l-1]}$$

$$A^{[l]} = \sigma(z^{[l]})$$

Conv (Input Data & rules)

DL (Input Data & Answers) expect ^{rules} Ans.

rules

Carlin

Page

/

/

For third layer

$$W^{[3]} = []_{1 \times 2}$$

$$Z^{[2]}$$

A^[2] → Finally is the

$$f(x, y, z) = x^2 + y^2 + z^2 - 8^2 = 0$$

Neuron that computes only linear thing, \rightarrow the one shared in Colab.

Other than that all neurons compute linear & non linear.

* Units :- no. of neuron in the layer

* input shape = [1]

Single neuron with multiple inputs performs similarly to logistic regression

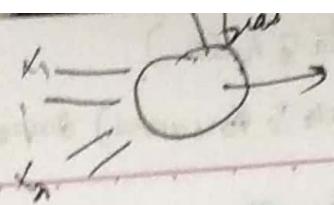
for eg: An image of a cat with 40K pixels.

Each pixel will be a separate input to a neuron

Now suppose we have 100K images each of 40K pixels.

actual data	I ₁	y
	I ₂	N
	I	Y
	I	Y
	I _{10K}	Y

There are 40K parameters (each pixel is param)



$$z^{(l)} = w^{(l)} a^{(l-1)} + b.$$

$$\hat{y} = f(z^{(l)})$$

assigning huge loss to

wrong prediction

i.e. upto ∞ .

$$L(y, \hat{y}) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

actual prob
(0 or 1)

consider $y=1, \hat{y}=1, L=0$

$y=0, \hat{y}=0, L=0$

y = actual or prediction value

\hat{y} = predicted value

Cost function = addition of all losses, then avg.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$$

where m is
number of losses

If all losses are 0, cost will also be 0.

Dimensions (parameters) in this case are 40×1 . (40 k are pixels & 6 is extra)

Gradient J'

$$w_i = w_i - \alpha \frac{\partial J(w, b)}{\partial w_i}$$

* is learning rate.

$$w_2 = w_2 - \alpha \frac{\partial J(w, b)}{\partial w_2}$$

positive slope so

go downward.

↓
↓

$$w_{40,000} = w_{40,000} - \alpha \frac{\partial J(w, b)}{\partial w_{40,000}}$$

$$b = b - \alpha \frac{\partial J(w, b)}{\partial b}$$

initially

We take ~~#~~ random values of w_1, w_2, \dots, w_n & then
 * for each parameter update it such that at last we
 have minimum value for each parameter.

* bias is some weightage.

which is not associated
 with the input.

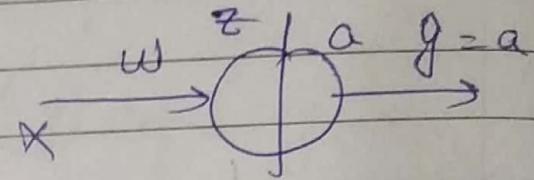
Vanilla version of gradient descent

25

30

$$L = -(y \log \hat{y} + (1-y) \log(1-\hat{y})),$$

$$w = w - \alpha \frac{dJ}{dw}$$



$$\frac{dL}{dw} = \frac{dL}{\hat{y}} \cdot \frac{d\hat{y}}{dz}$$

$$\frac{dL}{dw} = \frac{dL}{\hat{y}} \cdot \frac{d\hat{y}}{dz} \cdot \frac{dz}{dw}$$

$$\frac{dL}{\hat{y}} = \begin{pmatrix} y & + (1-y)(-1) \\ \hat{y} & (1-\hat{y}) \end{pmatrix}$$

$$= \frac{1-y}{1-\hat{y}} - \frac{y}{\hat{y}}$$

$$\hat{y} = \frac{1}{1+e^{-z}}$$

$$\frac{d\hat{y}}{dz} = \frac{e^{-z}}{(1+e^{-z})^2} = \hat{y}(1-\hat{y}).$$

$\frac{dz}{dw} = z$ since w is a single input so z is $-xw$.

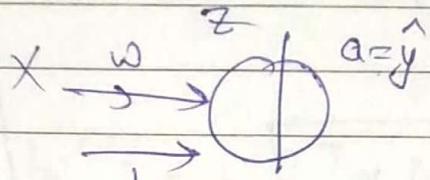
$$\frac{dz}{dw} = x.$$

$$\frac{dL}{dz} = \hat{y} - y_0 = (a - y)$$

$$\frac{dL}{dw} = \frac{dL}{dz} * \frac{dz}{dw}$$

$$= (a - y) * x$$

Back propagation \rightarrow Chain rule



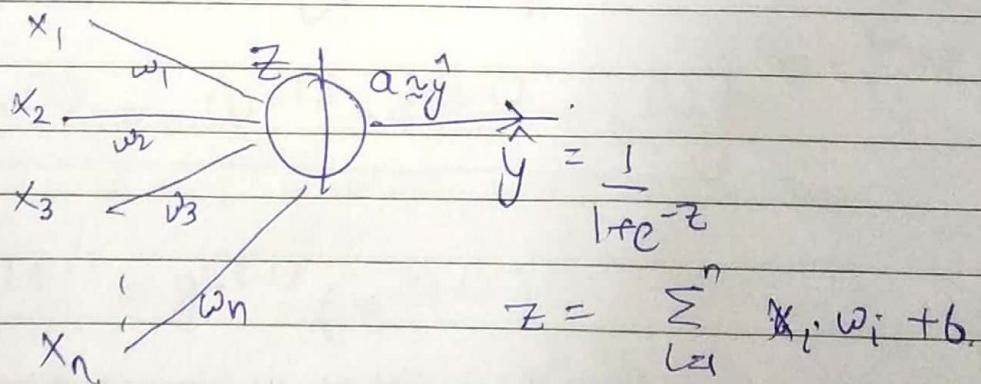
Now for f .

$$\frac{dL}{db} = \frac{dL}{dy} \frac{dy}{dz} \frac{dz}{db}$$

$$z = w \cdot x + b$$

$$\frac{dL}{db} = (\hat{y} - y)$$

* Back propagation logistic regression



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

$$\frac{dL}{dw_i} = \hat{y} - y \frac{dL}{dw_i}$$

$$z = x_1 w_1 + \dots + x_n w_n + b$$

$$\frac{dz}{dw_i} = x_i \frac{dx}{dw} \cdot \frac{dz}{dx_i}$$

$$\frac{dx}{dw} = (\hat{y} - y) \frac{dz}{dw}$$

↓
scalar

$$w = [w_1, w_2, \dots, w_n]^T$$

$$\frac{df}{dw} = (\hat{y} - y) x$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\boxed{\frac{dL}{dw} = (\hat{y} - y) x^T}$$

Forward & Backward

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$

Gradient Descent :-

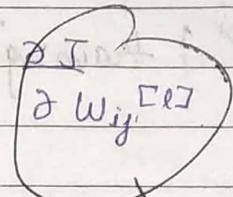
repeat → till loss change stops

for $l=1$ to L

for $i=1$ to $n^{[l]}$,

for $j=1$ to $n^{[l-1]}$

$$w_{ij}^{[l]} = w_{ij}^{[l]} - \alpha \frac{\partial J}{\partial w_{ij}^{[l]}}$$



$$b_i^{[l]} = b_i^{[l]} - \alpha \frac{\partial J}{\partial b_i^{[l]}}$$

→ obtained through back propagation.

Q Why we use gradient descent?

$$J = \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, \hat{y}^{(i)})$$

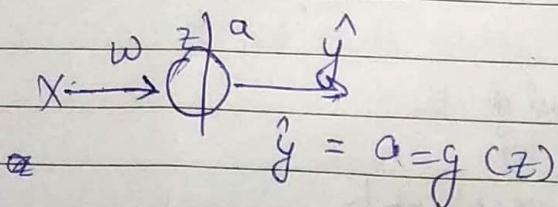
$$\hat{y} = a^{[L]} = g^{[L]}(z^{[L]}).$$

$$z^{[L]} = w^{[L]} \cdot a^{[L-1]} + b^{[L]}$$

activation function at layer L.

$$a^{[L-1]} = g^{[L-1]}(z^{[L-1]}).$$

$$z^{[L-1]} = w^{[L-1]} \cdot a^{[L-2]} + b^{[L-1]}$$



$$\hat{y} = a = g(z) = w \cdot x$$

$$\frac{dL}{d\hat{y}} \times \frac{dy}{dz} \times \frac{dz}{dw}$$

y

repeat

→ update converges

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

$$\alpha = (y, \hat{y}) = -(y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

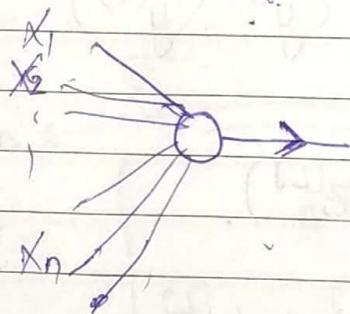
(y, \hat{y})

you don't have control over y
1 training example

$$y = e^{x^2+2}$$

$$\frac{dy}{dx} = e^{x^2+2} \cdot 2x$$

$$\frac{dL}{dw} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{dz} \times \frac{dz}{dw}$$



Now for more than 1 training example.

$w = \text{no. of neurons in } (dimensions) \times \text{no. of neurons in next layer} \times \text{no. of neurons in previous layer}$

Put T instead of t

$$z = \sum_{i=1}^n x_i w_i + b$$

Forward pass

$$\hat{y} = a = g(z)$$

$$z = \sum_{i=1}^n x_i w_i + b$$

From whom we get \hat{y} , ~~so~~ so loss is high we go for backward pass.

$$\left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n}, \frac{\partial L}{\partial b} \right]$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial z}{\partial w_1}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial z}{\partial w_2}$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial z}{\partial w_i}$$

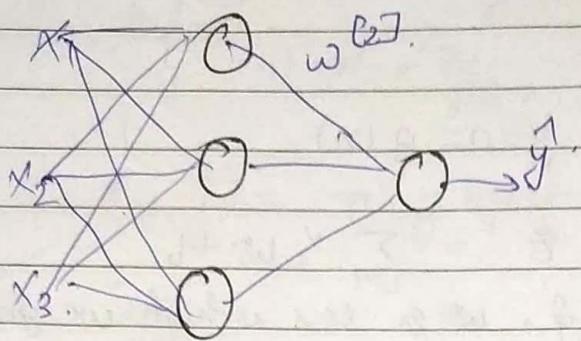
$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial z}{\partial b}$$

③ Vectorization

$$\frac{\partial L}{\partial w}[i] = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]$$

Instead of $w_i = w_i - \alpha \frac{\partial L}{\partial w_i}$

$$\text{write } w^{[l]} = w^{[l]} - \alpha \left[\frac{\partial L}{\partial w^{[l]}} \right]_{1 \times n}$$



$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_n^{[1]} \end{bmatrix} \quad n \times 1$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix}$$

$$\hat{y} = a^{[2]}$$

(Forward
Backward pass for
training example)

$$a^{[2]} = \begin{bmatrix} a_{11}^{[2]} \\ a_{12}^{[2]} \\ a_{13}^{[2]} \\ \vdots \\ a_{n1}^{[2]} \end{bmatrix}$$

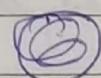
a.

$$z^{[i]} = w^{[i]} x + b^{[i]}$$

$$a^{[i]} = g(z^{[i]}),$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

~~step~~



$$z = \sum w_i x_i + b.$$

$$a = \sigma(z).$$

Why do we need backward pass:-

to calculate $\frac{\partial L}{\partial w_{ij}}$ or $\frac{\partial L}{\partial w_j}$.

Purpose of back propagation

Purpose is to tell network layer by layer from last to first \rightarrow that they are wrong.

Finally change your parameters occur directly.

$$w^T \frac{\partial L}{\partial w^{[i]}}$$

$$w^{[i]} = [w_1, w_2, \dots, w_n], \text{ (representation)}$$

$$\frac{\partial L}{\partial y} = \frac{(y - \hat{y})}{\hat{y}(1-\hat{y})}$$

in matrix form.

$$\frac{\partial L}{\partial \hat{y}} = \frac{1-\hat{y}}{1-\hat{y}} - \frac{\hat{y}}{\hat{y}}$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{d}{dz} \left(\frac{1}{1+e^{-z}} \right) = \frac{-1 \times e^{-z}}{(1+e^{-z})^2}$$

$$\frac{-1 \times e^{-z}}{(1+e^{-z})^2} = \frac{1}{1+e^{-z}} - 1$$

$$\frac{\partial y}{\partial z} = \left(\frac{1}{1+e^{-z}} \right) * \frac{(1+e^{-z}-1)}{(1+e^{-z})}$$

$$= g(z) (1-g(z))$$

5
 $\frac{\partial y}{\partial w_i}$

10
 $\frac{\partial z}{\partial w_i} = x_i$

we clubbed it because .

$$\Rightarrow \frac{\partial L}{\partial w_i} = \left(\frac{\partial L}{\partial z} \right) x_i$$

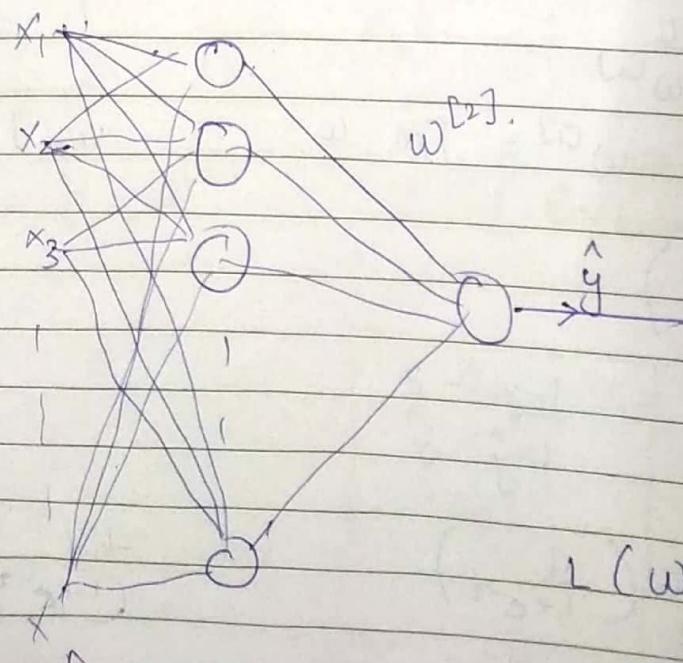
trying to use matrices (vectorizing).

15
 Ans:

$$\frac{\partial L}{\partial w^{[i]}} = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]$$

20

$$\left[\frac{\partial L}{\partial z} x_1, \frac{\partial L}{\partial z} x_2, \dots, \frac{\partial L}{\partial z} x_n \right]$$



30
 Fig 1.

$$L(w^{[2]}, b^{[2]}, w^{[1]}, b^{[1]})$$

$$\frac{\partial L}{\partial \hat{y}} = \frac{\hat{y} - y}{(1-\hat{y}) \hat{y}}$$

Now we want $\left(\frac{\partial L}{\partial g}\right) \left(\frac{\partial g}{\partial z}\right)$

$$= \frac{\hat{y} - y}{(1-\hat{y}) \hat{y}} * (\hat{y}(1-\hat{y})) = \hat{y} - y$$

$$\boxed{\frac{\partial L}{\partial w} = (\hat{y} - y) x^T}$$

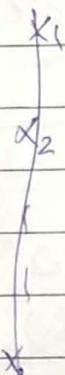
$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

logistic regression

dimension of w is $1 \times n$.

so we have to take x^T to multiply with $\hat{y} - y$.

$$\boxed{\frac{\partial L}{\partial b} = \hat{y} - y}$$



Forward propagation for Fig 1

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

$$z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}$$

we can use different

$$a^{[2]} = g^{[2]}(z^{[2]})$$

activation functions
at diff layers

$$y = y^b$$

Backward propagation

$$\frac{\partial L}{\partial w^{[2]}} = \left(\frac{\partial L}{\partial z^{[2]}} \right) \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

$$= (y - y)$$

$$= (a^{[2]} - y) \cdot a^{[2] T}$$

$$\frac{\partial L}{\partial b^{[2]}} = a^{[2]} - y$$

$$\frac{\partial L}{\partial w^{[1]}}$$

$$w_i = w_i - \frac{\partial L}{\partial w^{[1]}}$$

$$\frac{\partial L}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial b^{[2]}} = a^{[2]} - y$$

Now moving a

step back

$$\frac{\partial L}{\partial w^{[1]}} = \frac{\partial L}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial w^{[1]}}$$

$$\frac{\partial L}{\partial z^{[1]}} = w^{[2] T} \left(\frac{\partial L}{\partial z^{[2]}} \right) \cdot g^{[1]}(z^{[1]})$$

looking
for its
derivative

Derivative

$$\frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

$$\frac{\partial L}{\partial z^{[1]}} = \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}}$$

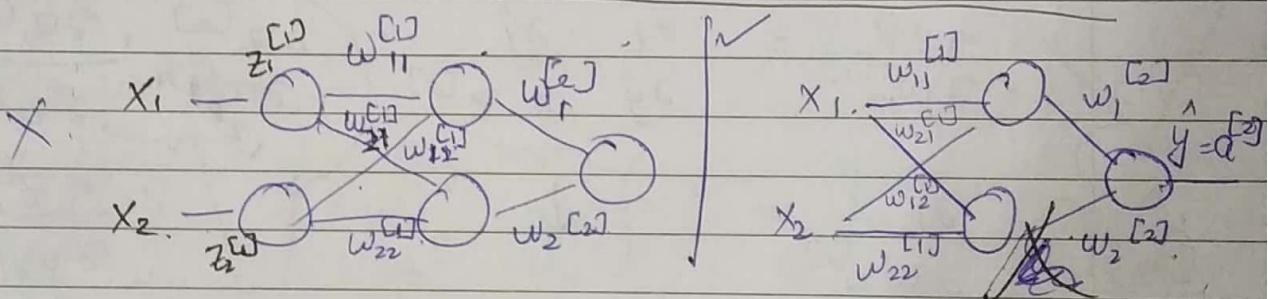
$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial w_{12}}$$

Summary of gradient descent

Now has to derive this

$$\cancel{w^{[2]}} = w^{[2]} = \frac{\partial L}{\partial z^{[2]}} \begin{bmatrix} w_1^{[2]} \\ w_2^{[2]} \\ \vdots \\ w_n^{[2]} \end{bmatrix} \times \begin{bmatrix} g_1^{[1]} \\ z_1^{[1]} \\ \vdots \\ g_n^{[1]} \end{bmatrix}$$

$$w^{[2]} = [w_1^{[2]}, w_2^{[2]}, \dots, w_n^{[2]}]$$



$$\frac{\partial L}{\partial w^{[2]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial w^{[2]}}$$

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

$$\hat{y} = a^{[2]} = g(z^{[2]}).$$

$$\frac{\partial L}{\partial w^{[1]}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \cdot \frac{\partial a^{[1]}}{\partial z^{[1]}} \cdot \frac{\partial z^{[1]}}{\partial w^{[1]}}$$

This is final ans :-

$$d\hat{z}^{[1]} = w^{[2] \top} d\hat{z}^{[2]} \rightarrow g^{[1]'}(z^{[1]})$$

$$\frac{\partial L}{\partial w_1^{[2]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial w_1^{[2]}}$$

~~$$\frac{\partial L}{\partial w_2^{[2]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial w_2^{[2]}} * \frac{\partial z^{[1]}}{\partial z^{[2]}} * \frac{\partial z^{[1]}}{\partial w_2^{[1]}}$$~~

$$\frac{\partial L}{\partial w_3^{[2]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial w_3^{[2]}}$$

$$\frac{\partial L}{\partial w_{11}^{[1]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial a_1^{[0]}} * \frac{\partial a_1^{[0]}}{\partial z_1^{[1]}} * \frac{\partial z_1^{[1]}}{\partial w_{11}^{[0]}}$$

$$\frac{\partial L}{\partial w_{12}^{[1]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial a_1^{[0]}} * \frac{\partial a_1^{[0]}}{\partial z_2^{[1]}} * \frac{\partial z_2^{[1]}}{\partial w_{12}^{[0]}}$$

$$\frac{\partial L}{\partial w_{21}^{[1]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial a_2^{[0]}} * \frac{\partial a_2^{[0]}}{\partial z_1^{[1]}} * \frac{\partial z_1^{[1]}}{\partial w_{21}^{[0]}}$$

$$\frac{\partial L}{\partial w_{22}^{[1]}} = \frac{\partial L}{\partial y} * \frac{\partial y^1}{\partial z^{[2]}} * \frac{\partial z^{[2]}}{\partial a_2^{[0]}} * \frac{\partial a_2^{[0]}}{\partial z_2^{[1]}} * \frac{\partial z_2^{[1]}}{\partial w_{22}^{[0]}}$$

Now we want final ans a

$$d_{z^{[1]}} = w^{[2]^T} * d_{z^{[2]}} * g^{[0]} * (z^{[0]})$$

$$z^{[2]} = w^{[2]} a^{[0]} + b^{[2]}$$

$$= [w_1^{[2]} \quad w_2^{[2]}] \begin{bmatrix} a_1^{[0]} \\ a_2^{[0]} \end{bmatrix} + b^{[2]}$$

element
wie

$$w_1^{[2]} \cdot a_1^{[1]} + w_2^{[2]} \cdot a_2^{[1]}$$

$$\frac{\partial z^{[2]}}{\partial a_1^{[1]}} = w_1^{[2]}$$

$$\frac{\partial z^{[2]}}{\partial a_2^{[1]}} = w_2^{[2]}$$

This
is the idea
where the terms can
come from.

$$\left. \begin{aligned} dz^{[1]} &= w^{[2]T} dz^{[2]} * g'(z^{[1]}) \\ \frac{\partial z^{[2]}}{\partial a_1^{[1]}} & \end{aligned} \right\}$$

Now we will look exactly where
these terms comes from

$$z^{[1]} = w^{[1]T} x + b^{[1]}$$

we

are finding

$$\frac{\partial L}{\partial z^{[1]}} \left[\begin{array}{c} z_1^{[1]} \\ z_2^{[1]} \end{array} \right] = \left[\begin{array}{cc} w_{11}^{[1]} & w_{12}^{[1]} \\ w_{21}^{[1]} & w_{22}^{[1]} \end{array} \right] \left[\begin{array}{c} x_1 \\ x_2 \end{array} \right] + \left[\begin{array}{c} b_1^{[1]} \\ b_2^{[1]} \end{array} \right] \frac{\partial a_1^{[1]}}{\partial z^{[1]}}$$

$$\left. \begin{aligned} \left[\begin{array}{c} z_1^{[1]} \\ z_2^{[1]} \end{array} \right] &= w_{11}^{[1]} x_1 + w_{12}^{[1]} x_2 + b_1^{[1]} \\ \left[\begin{array}{c} z_1^{[1]} \\ z_2^{[1]} \end{array} \right] &= w_{21}^{[1]} x_1 + w_{22}^{[1]} x_2 + b_2^{[1]} \end{aligned} \right\}$$

element

wise

$$\frac{\partial L}{\partial z_1^{[1]}} = dz_1^{[1]} = (a^{[2]} - y) \cdot w_1^{[2]} \cdot g'(z_1^{[1]})$$

$$\left. \begin{aligned} \left[\begin{array}{c} \frac{\partial L}{\partial z_1^{[1]}} \\ \frac{\partial L}{\partial z_2^{[1]}} \end{array} \right] &= \left[\begin{array}{c} (a^{[2]} - y) \cdot w_1^{[2]} \cdot g'(z_1^{[1]}) \\ (a^{[2]} - y) \cdot w_2^{[2]} \cdot g'(z_2^{[1]}) \end{array} \right] \end{aligned} \right\}$$

Try this again

$$\sigma(n) = \frac{1}{1+e^{-n}}$$

$$= (a^{(2)} - y) \left[w_1^{(2)} \cdot g_1^{(1)} (z_1^{(1)}) + w_2^{(2)} \cdot g_2^{(1)} (z_2^{(1)}) \right]$$

$$= (a^{(2)} - y) \begin{bmatrix} w_1^{(2)} \\ w_2^{(2)} \end{bmatrix} \times \begin{bmatrix} g_1'(z) \\ g_2'(z) \end{bmatrix}$$

$$w^{(2)} = \begin{bmatrix} w_1^{(2)} & w_2^{(2)} \end{bmatrix}$$

$$w^{(2)T} = \begin{bmatrix} w_1^{(2)} \\ w_2^{(2)} \end{bmatrix}$$

Back propagation \rightarrow if any of the derivative is small.

$$w_{ij}^{(l)} = w_{ij}^{(l)} - \alpha \frac{\partial J}{\partial w_{ij}^{(l)}}$$

- \rightarrow Vanishing gradient problem \rightarrow No soln
- \rightarrow exploding gradient problem

Solution to this is Gradient clipping

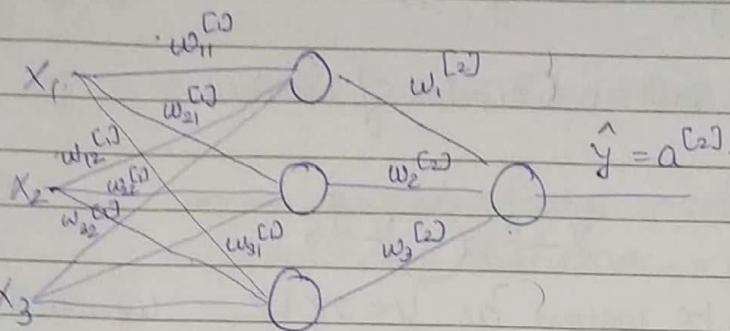
If gradient comes close to 0, it might not be the optimal soln but your learning will stop.

H.W.

3 nodes in first layer
1 node in last layer

$$Z = wX + b \quad \frac{\partial L}{\partial b} = \frac{\partial L}{\partial Z}$$

Camlin Page
Date / / 1
3 inputs



~~$\frac{\partial L}{\partial a^{[2]}}$~~

$$\begin{bmatrix} \frac{\partial L}{\partial a^{[2]}} \\ \dots \\ \frac{\partial L}{\partial a^{[2]}} \end{bmatrix}$$

1×3

$$\begin{bmatrix} \frac{\partial L}{\partial z^{[2]}} \\ \dots \\ \frac{\partial L}{\partial z^{[2]}} \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} w_1^{[2]} & w_2^{[2]} & w_3^{[2]} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial L}{\partial z^{[2]}} \\ \dots \\ \frac{\partial L}{\partial z^{[2]}} \end{bmatrix}$$

$$w^{[2]T} = \begin{bmatrix} w_1^{[2]} \\ w_2^{[2]} \\ w_3^{[2]} \end{bmatrix} \begin{bmatrix} g'(z_1^{[2]}) \\ g'(z_2^{[2]}) \\ g'(z_3^{[2]}) \end{bmatrix}$$

$$\frac{\partial L}{\partial w_{11}^{[1]}}$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial a_1^{[1]}} \cdot \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \cdot \frac{\partial z_1^{[1]}}{\partial w_{11}^{[1]}}$$

$$\frac{\partial L}{\partial w_{21}^{[1]}}$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial a_2^{[1]}} \cdot \frac{\partial a_2^{[1]}}{\partial z_2^{[1]}} \cdot \frac{\partial z_2^{[1]}}{\partial w_{21}^{[1]}}$$

$$\frac{\partial L}{\partial w_{12}^{[1]}}$$

$$= \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z^{[2]}} \cdot \frac{\partial Z^{[2]}}{\partial a_1^{[1]}} \cdot \frac{\partial a_1^{[1]}}{\partial z_1^{[1]}} \cdot \frac{\partial z_1^{[1]}}{\partial w_{12}^{[1]}}$$

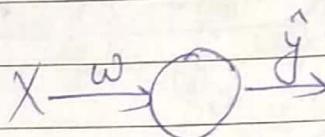
$$(a^{[2]}\hat{y}) * g'(z_1^{[1]})$$

$$Z^{[2]} = w^{[2]} \cdot X + b^{[2]}$$

Code

nof (keras.layers.Dense) are no of layers.

units=1 (number of neurons).



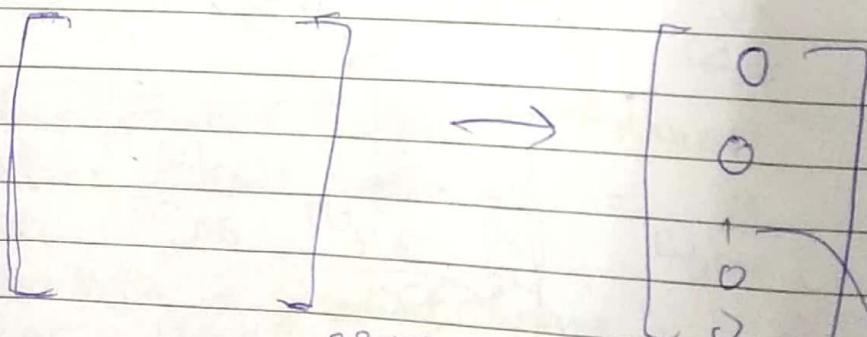
f^n can be written as $y = 2x + 2$. We want our neural network to learn the function.

model.compile (optimizer='SGD', loss='mean_squared_error')

model.fit (fit giving input & output trying to learn the fn,

Fashion MNIST data

$$\hat{y} = f(x)$$



images a dimension

denoted in
form of one
hot vector

Dense suggest that every neuron in the layer have inputs from all features.

* keras.layers.Dense means all inputs are connected to this layer (each neuron)

* 90% of time we use relu (activation function) \rightarrow hidden layer

loss fn

\leftarrow binary classification \rightarrow softmax or sigmoid
multiclass classification \rightarrow softmax

* tf.keras.models.Sequential ()

all these layers are executed sequentially

Q Why can't we use softmax as activation fn in hidden layer.

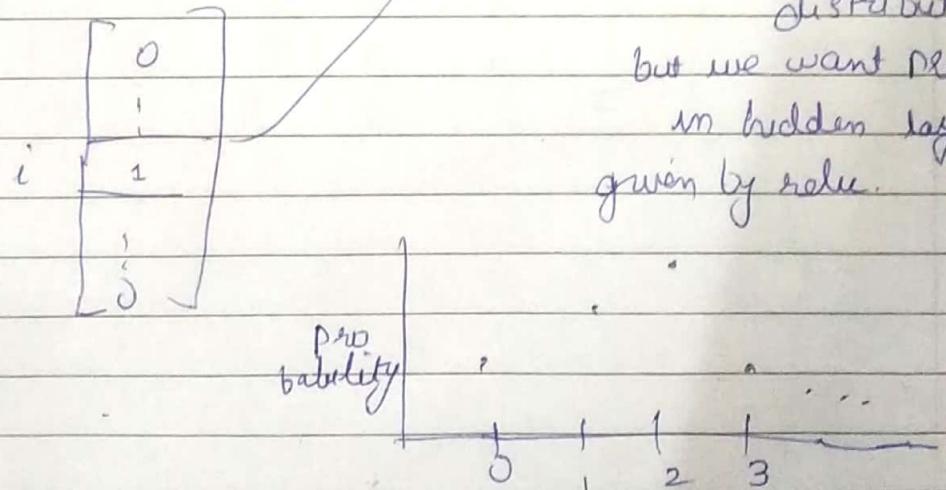
We can use but problem is .

tf. keras

softmax looks like

$$\frac{e^{x_i}}{\sum e^{x_i}}$$

e^x \rightarrow will give probability distribution



but we want probability in hidden layer given by relu.

1 epoch :- do it once for whole data set , it is 1 epoch.

Initially there was a while loop that runs until it updates all. To achieve this, we use epochs which ensures that no. we will run till epochs (whatever the no be)

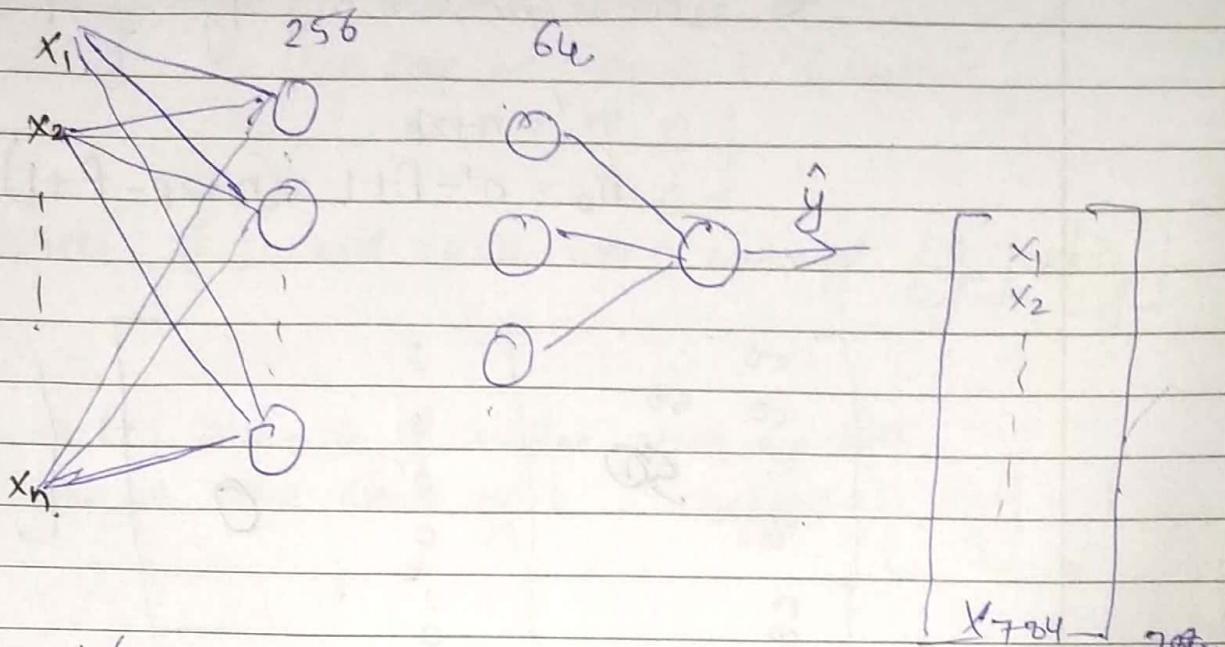
- * Solve exercises from Jupyter
- * Model starts learning from where it was left.

what is use of CNN?

Nothing special but will be able to do same thing with less parameters.

One thing CNN do is \rightarrow don't flatten the image like that of deep neural networks.

tf.keras.dense.flatten \rightarrow what it does is take the 28×28 image & flatten it.



→ keeps/preserves local features in an image.

→ less # parameters.

$$\begin{array}{c}
 \begin{array}{c} 5 & 3 & 6 & 8 & 9 & 1 & 7 \\ 1 & 6 & 8 & 3 & 2 & 1 & 1 \\ 9 & 9 & 3 & 2 & 1 & 1 & 1 \\ 1 & 1 & 8 & 4 & 6 & 5 & 1 \\ 9 & 3 & 4 & 0 & 0 & 1 & 1 \\ 0 & 1 & 2 & 3 & 4 & 5 & 1 \end{array} \\
 \times \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & 1 \\ 1 & 0 & -1 \end{bmatrix} \\
 = \begin{bmatrix} -2 & 5 & 5 & 10 \\ -8 & 7 & & \\ -4 & & -4 \end{bmatrix} \quad 6 \times 6
 \end{array}$$

$n \times n$ matrix with a $f \times f$ filter $\rightarrow (n-f+1) \times (n-f+1)$.

If we want to increase contribution from 5, we pad 0's everywhere.

Padding P will result into new dimension $n' = n + 2P$

Vertical filter lets you

recognize vertical edges

→ Solve exercises on filter if any.

$$n' = n + 2p$$

$$n_o = n' - f + 1 = (n + 2p - f + 1) \times (n + 2p - f + 1)$$

$$\left[\begin{array}{cc|c} 50 & & 0 \\ 50 & 50 & 0 \\ 50 & & 0 \\ 50 & & 0 \\ \hline 50 & & 0 \end{array} \right] \quad \left[\begin{array}{ccc} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{array} \right]$$

if transition → $\left[\begin{array}{cccc} 0 & 150 & 150 & 0 \\ 0 & 150 & 150 & 0 \\ 0 & & & 0 \\ 0 & & & 0 \end{array} \right]$

from light to dark

if transition from dark to light

$$\left[\begin{array}{cccc} 0 & -150 & -150 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{array} \right]$$

$$\left[\begin{array}{cccc} 0 & 150 & 150 & 0 \\ 0 & 150 & 150 & 0 \\ 0 & 150 & 150 & 0 \\ 0 & 150 & 150 & 0 \end{array} \right] \left[\begin{array}{cc} 1 & -1 \\ -1 & -1 \end{array} \right] = \left[\begin{array}{ccc} -300 & 0 & 300 \\ -300 & 0 & 300 \\ -300 & 0 & 300 \\ -300 & 0 & 300 \end{array} \right]$$

4x4 3x3

But suppose if we want to study

$$\left\lfloor \frac{n-f+1}{s} \right\rfloor \quad \left\lfloor \frac{n-f}{s} \right\rfloor + 1$$

$$\left\lfloor \frac{n-f+1}{s} \right\rfloor$$

In that case we want 2×2 matrix

$$\left[\begin{smallmatrix} 1 & -1 \\ 1 & -1 \end{smallmatrix} \right]$$

Stride of 2 will result into a matrix of $\left[\begin{smallmatrix} 1 & -1 \\ 1 & -1 \end{smallmatrix} \right] \times \left[\begin{smallmatrix} n-f+1 \\ s \end{smallmatrix} \right]$

$n-f+1$ was size of matrix after filter.

now we want stride of 2. ~~n-f+1~~ $\left[\begin{smallmatrix} n-f+1 \\ s \end{smallmatrix} \right]$

so if we want stride as :

$$\left[\begin{array}{cccc} 0 & 180 & 180 & 0 \\ 0 & 180 & 180 & 0 \\ 0 & 180 & 180 & 0 \\ 0 & 180 & 180 & 0 \end{array} \right] \left[\begin{array}{c} 1 \\ -1 \\ 1 \\ -1 \end{array} \right].$$

$$= \left[\begin{array}{c} \square \\ \square \end{array} \right]_{2 \times 2}$$

Stride is should be less than filter.

$$\begin{matrix} 0 & \boxed{180} & 180 & 0 \\ 0 & 180 & 180 & 0 \\ 0 & 180 & 180 & 0 \\ 0 & 180 & 180 & 0 \end{matrix}$$

$\uparrow \leftrightarrow \uparrow$ *batch.*
 $0 \quad 0 \quad 0$ *(cannot be distinguished)*

So if size is less than overlapping will be there

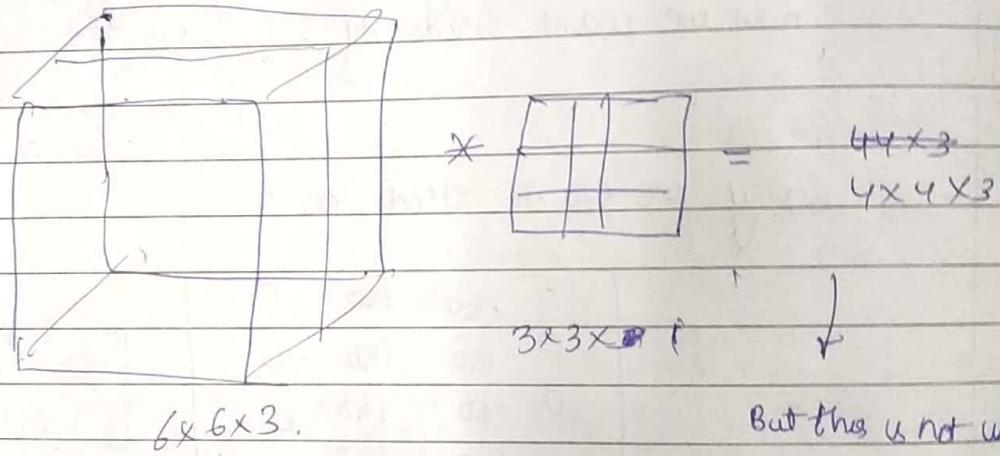
Convolution in colored image

$(R, G, B) \rightarrow 3 \text{ bytes}$

$n \times n \times 1 \rightarrow \text{black \& white image}$

$n \times n \times 3 \rightarrow \text{colored image}$

Q How to take convolution for 'colored image'?



But this is not used.

$$6 \times 6 \times 3. \quad 3 \times 3 \times 3. = 4 \times 4 \times 1.$$

↓
This is used.

$n_0^{[C_l]} \times n_1^{[C_l]} \times n_2^{[C_l]}$ → channel (e, g, b)
 $p^{[C_l]}, s^{[C_l]}, f^{[C_l]}$
padding stride filter size in
in layer l in layer l layer l

$$\begin{matrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{matrix}$$

$$n_w^{[0]} \times n_h^{[0]} \times n_c^{[0]}$$

↓ ↓ ↓

f 6 3

After convolution :-

$$n_w^{[1]} \times n_h^{[1]} \times n_c^{[1]}$$

↓ ↓ ↓

4 4 1

$$P^{[0]} = 0$$

$$S^{[0]} = 1$$

$$f^{[1]} = 3 \times 3 \times n_c^{[1]}$$

Why 3×3 ? \rightarrow It is nice to have a central element & all rest around it.

$$\begin{array}{l} n_w \times n_h \times n_c \\ \uparrow \\ 6 \times 6 \times 3 \end{array} \times f_1 \rightarrow 3 \times 3 \times 3 (f_1) \rightarrow 4 \times 4 \times 1$$

$$\times f_2 \rightarrow 4 \times 4 \times 1$$

$$\times f_3 \rightarrow 4 \times 4 \times 1$$

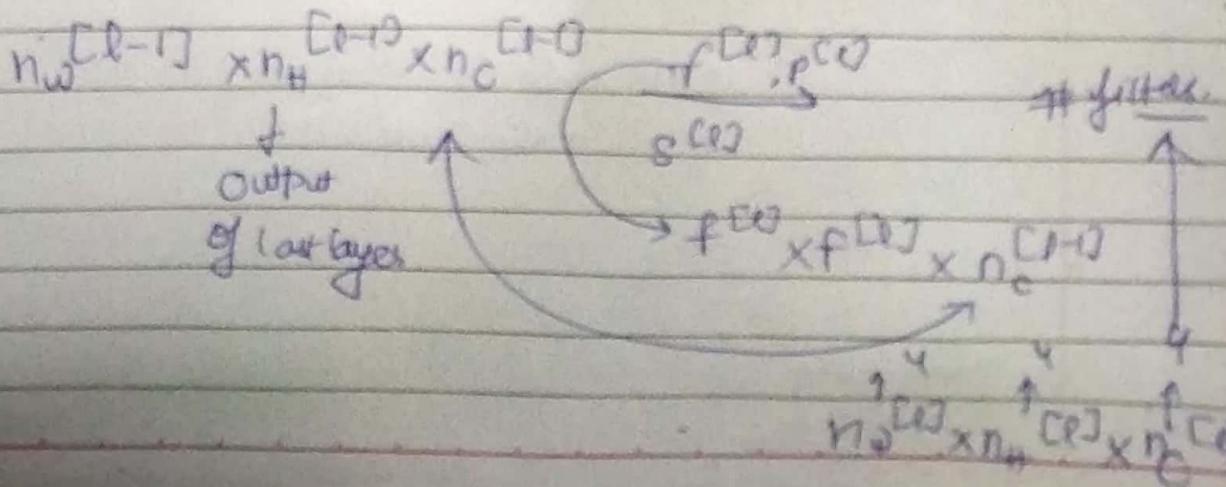
$$\times f_4 \rightarrow 4 \times 4 \times 1$$

These gives different aspects
of the image

\rightarrow for book keeping stack
them over others.

$$4 \times 4 \times 4$$

$$n_w^{[0]} n_h^{[0]} n_c^{[0]}$$



$$b^{[e]} = \begin{cases} \# filters \end{cases}$$

linear pass
is done for
conv.

$$z^{[e]} = w^{[e]} A^{[e-1]}$$

$$+ b^{[e]}$$

We want
to map it here

$$A^{[e]} = g^{[e]}(z^{[e]})$$

parameters are filters

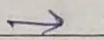
take biases by adding.

Bias does not depend on filters.

$$w_1 \quad w_2 \quad w_3$$

$$w_4 \quad w_5 \quad w_6$$

$$w_7 \quad w_8 \quad w_9$$



We have to

have non linear

If we don't
have it

$$n_w^{[l-1]} \times n_h^{[l-1]} \times n_c^{[l-1]}$$

$$\rightarrow n_w^{[l]} \times n_h^{[l]} \times n_c^{[l]}$$

$$w = \sigma(h(g(j(x))))$$

We need non linear

in every layer because

$$\text{Relu}(n_w^{[e]} \times n_h^{[e]} \times n_c^{[e]})$$

$y < 0$ make 0

$y > 0$ make same

Every time in each layer we use filters (parameters).

CNN to DNN.

Pooling ↗

$$\begin{bmatrix} 6 & 4 & 2 & 1 \\ 9 & 3 & 2 & 8 \\ 1 & 8 & 7 & 1 \\ 2 & 0 & 3 & 7 \end{bmatrix} = \begin{bmatrix} 9 & 6 \\ 5 & 7 \end{bmatrix}_{2 \times 2}$$

training
validation
test Improve accuracy
of all.

We do max pooling because we want to extract bright side of an image.

Pooling preserve translation to some extent. → Translation invariant.
" does not preserve rotation.

→ you can make pooling learn the invariant.

In pooling: we keep a 2d filter (imagine). $f \times f$.

There are no learning parameters.

$3 \times 3 \times 3$

10 filters ~~or 100~~

$27 \times 10 + 10 = 280$ parameters.

Parameters that you want to learn are invariant to inputs so even if $200 \times 200 \times 3$ image, we still have 280 parameters.

But computational cost increases.

But in Deep neural network ~~the~~ parameters are not invariant of input size, that's why we use CNN.

↓

