

What is Deep Learning

Deep learning is a part of machine learning that uses artificial neural networks with multiple layers to learn patterns from data.

Key idea

Instead of humans designing features manually, deep learning model automatically discover the best features during training. This makes them powerful for solving non linear data as well.

Importance

- ① Powers modern AI applications
- ② Works well with large dataset
- ③ Improve with better hardware like GPU's.

Applications

- ① Image recognition
- ② Natural language processing.
- ③ Speech recognition
- ④ Recommendation systems etc.

Difference between ML and DL.

practical and useful

ML

- ① Works well with smaller dataset
- ② Features are manually design by humans.
- ③ Use methods like decision tree, SVM, regression etc.
- ④ Requires less computational power
- ⑤ Great explainability

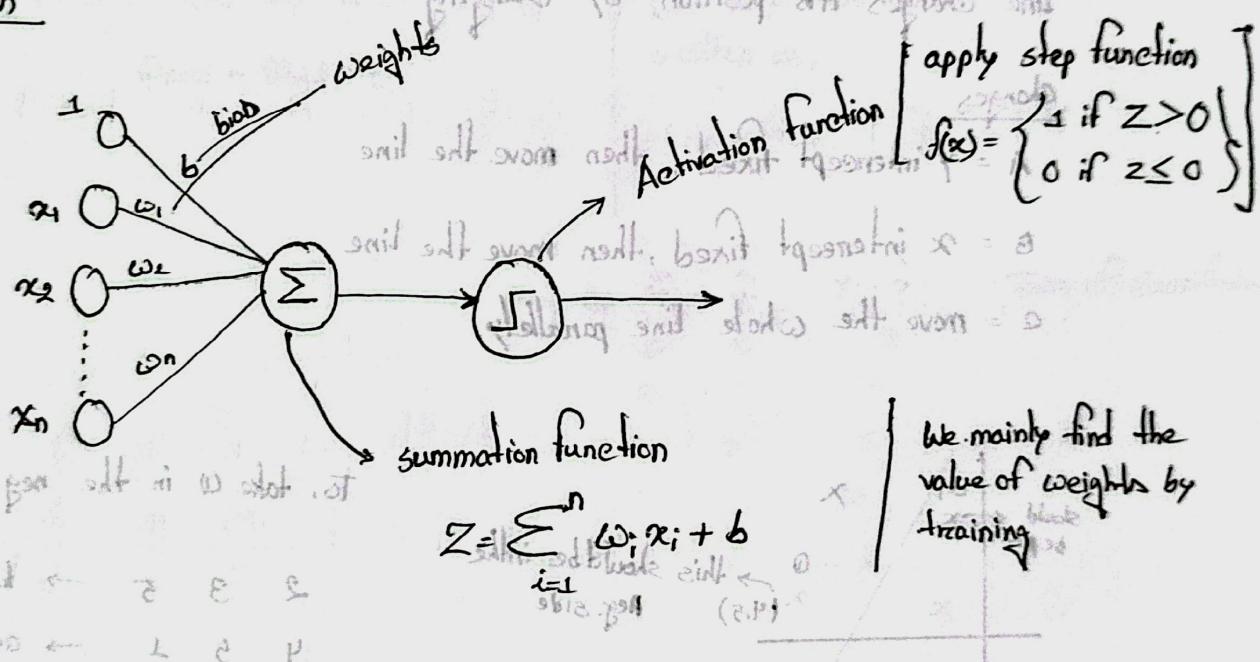
DL

- ① Needs large amount of data.
- ② Features are learned automatically by models
- ③ Uses multi-layer neural networks.
- ④ Needs GPUs / TPUs for training
- ⑤ Less explainability

④ Perception

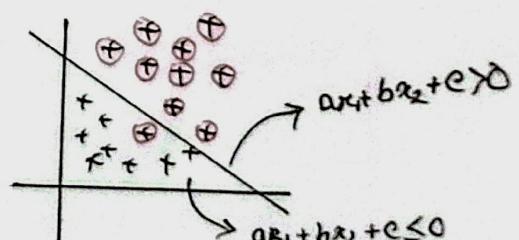
- ⇒ Simplest artificial neural networks.
- ⇒ It's a binary classifier.
- ⇒ Learns using a supervised learning rules.
- ⇒ Works if a linear boundary exists (AND, OR)
- ⇒ Cannot solve problems that are not linearly separable. (XOR)

Diagram



For 2D, it will create a line which separates the classes.

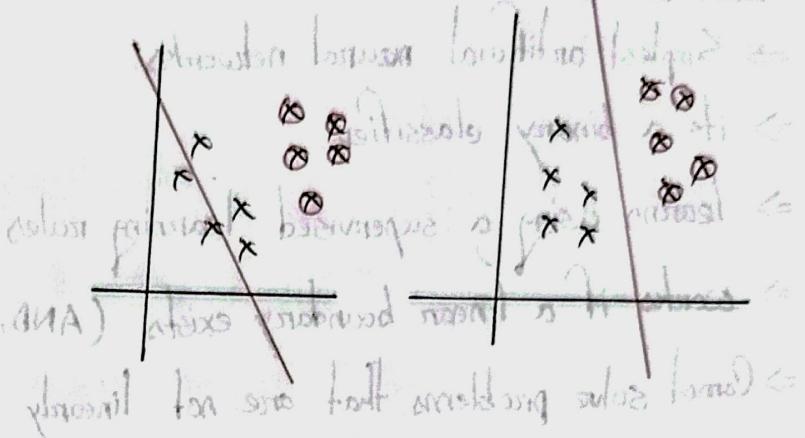
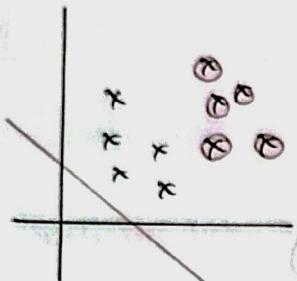
For 3D, → Plane
Above, $P \rightarrow$ hyperplane



For these kind of graph, not ideal, so we send it out to expand off.
Not linearly separable.



How it creates the line



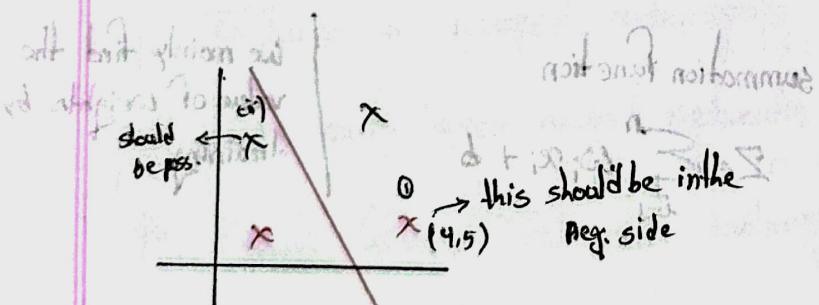
Line changes its position by changing the value of A, B, and C.

Changes

A = y intercept fixed, then move the line

B = x intercept fixed, then move the line

C = move the whole line parallelly,



To, take C in the neg.

2 3 5 → line coeff.

4 5 1 → coordinate and extra 1

$$\frac{2x + 3y + 5 = 0}{-2x - 2y - 2 = 0}$$

new line 19 → , 40, 19

$$-2x - 2y + 4 = 0$$

For, (ii) add, instead of subtract.

The changes of line is huge, so we use learning rate (4-7)

, step size of gradient fall

Algorithm

epoch \rightarrow 1000, $\eta = 0.01$

for $i \in \text{range}(\text{epoch})$:

take a random row data

if $(x_i \in N \text{ and } z \geq 0) \vee (x_i \in P \text{ and } z \leq 0)$

$$w_{\text{new}} = w_{\text{old}} - \eta x_i$$

}

elif $(x_i \in P \text{ and } z \leq 0) \vee (x_i \in N \text{ and } z \geq 0)$

$$w_{\text{new}} = w_{\text{old}} + \eta x_i$$

}

(cost) \rightarrow weight cost

$$((w_0 + w_1 x_1 + w_2 x_2) x_{\text{out}}) \sum_{i=1}^n \frac{1}{n} = J$$

$$w = [w_0, w_1, w_2]$$

A A B
B A C

x_i = class

η = learning rate

Total if and elif part can be written as,

$$w_n = w_0 + \eta (y_i - \hat{y}_i) x_i$$

$$\frac{y_i}{0} \quad \frac{\hat{y}_i}{0} \quad \frac{y_i - \hat{y}_i}{0} \quad \text{Term of classification}$$

$$\begin{matrix} 0 \\ 0 \\ 1 \\ -1 \\ 1 \\ 1 \\ 0 \end{matrix}$$

$$[(w_0 + w_1 x_1 + w_2 x_2) x_{\text{out}}] \sum_{i=1}^n \frac{1}{n} = J$$

$\frac{y_i}{0}$	$\frac{\hat{y}_i}{0}$	$\frac{y_i - \hat{y}_i}{0}$
0	0	0
0	1	-1
1	1	0
1	0	1
0	0	0
0	1	-1
1	0	1
1	1	0

$$J = \text{sum} = (2 \cdot 0, 0) x_{\text{out}} \quad ①$$

$$J = \text{sum} = (2 \cdot 1, 0) x_{\text{out}} \quad ②$$

$$J = \text{sum} = (2 \cdot 1, 1) x_{\text{out}} \quad ③$$

$$J = \text{sum} = (2 \cdot 1, 1) x_{\text{out}} \quad ④$$

0 = student has contributing features, not sufficient features gets no

Loss function (Hinge)

$$L = \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

We need to find such a line which loss function value will be minimum.

$$L = \underset{\omega_0, \omega_1, \omega_2, b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, y_i f(x_i))$$

To do that, we will use gradient descent.

Explanation of Loss function

for, 2 sample, the loss function will look like,

$$L = \frac{1}{2} \left[\max(0, -y_1 f(x_1)) + \max(0, -y_2 f(x_2)) \right]$$

$$\textcircled{1} \quad \max(0, -1 \cdot 1) = -\text{ve} = 0$$

$$\textcircled{2} \quad \max(0, -1 \cdot (-1)) = +\text{ve} = 0$$

$$\textcircled{3} \quad \max(0, +1 \cdot 1) = +\text{ve} = 0$$

$$\textcircled{4} \quad \max(0, (+1 \cdot -1)) = +\text{ve} = 0$$

For, correct prediction, contribute = 0
or else, atleast something.

10.1	$n = \text{number of sample}$ $y_i = \text{actual output}$ $f(x_i) = z = \text{predicted output}$ $z = \omega_0 + \omega_1 x_1 + \omega_2 x_2$
------	---

Data		
x_1	x_2	y
x_{11}	x_{12}	y_1
x_{21}	x_{22}	y_2

y_i	\hat{y}_i	$\boxed{0 = -1}$
1	1	① \rightarrow correct
1	-1	②
-1	1	③
-1	-1	④ \rightarrow correct

Benefits

Perceptron is very flexible. We can change or modify loss function and activation function for different purpose.

Here are some of the combination

Loss function

Hinge Loss

Log Loss (Binary cross entropy)

Categorical Cross entropy

Mean Squared error (MSE)

Activation

Step activation

SIGMOID

Softmax

Linear (nothing)

output (works like)

perceptron \rightarrow binary classi.

Logistic regression (Binary)

Softmax regression (multi classi.)

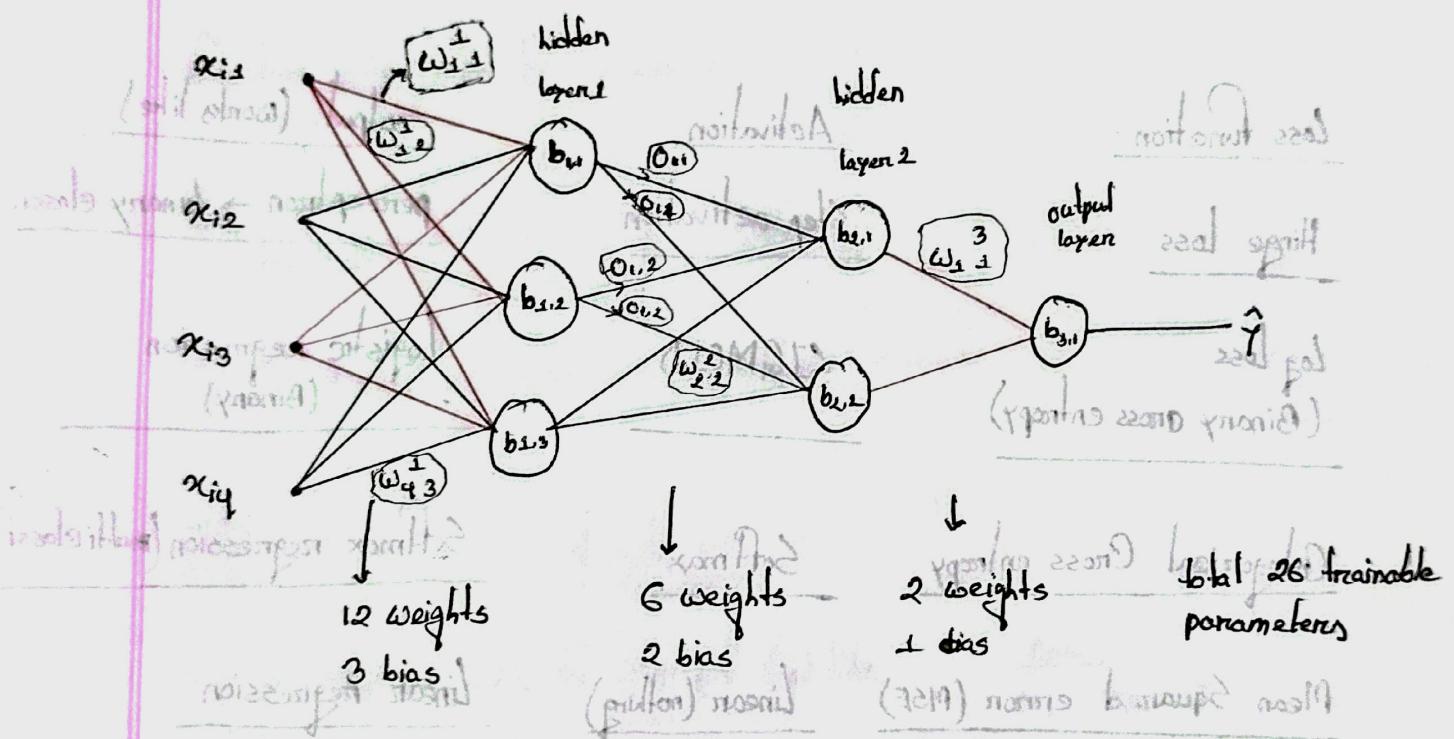
Linear regression

- linear function \rightarrow
- step function \rightarrow
- sigmoid function \rightarrow

$x \in \mathbb{R} \leftarrow$ input

Multi Layered perceptron (MLP) Notation

Notation used, which no specific one is standard, this is consequent
 Inputs x_{ij} (real numbers)
 Weights w_{ij}^k (real numbers)
 Bias b_{ij} (real numbers)
 Output o_{ij} (real numbers)



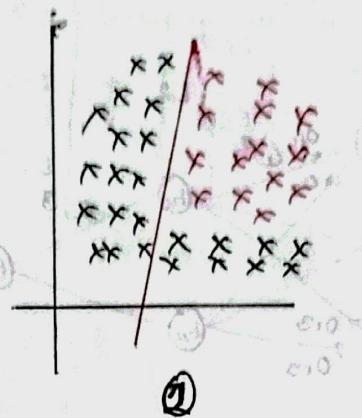
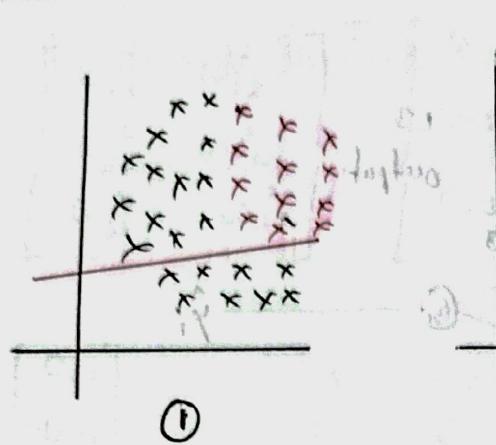
Notation

Bias $\Rightarrow \underline{b_{ij}}$ $i = \text{layer no.}$
 $j = \text{node no.}$

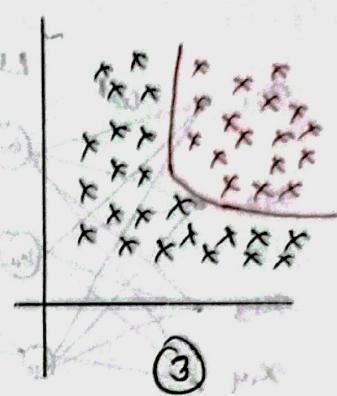
Output $\Rightarrow \underline{o_{ij}}$

Weights $\Rightarrow \underline{\underline{w_{ij}^k}}$ $k = \text{entering layer no.}$
 $i = \text{previous layer node no.}$
 $j = \text{entering layer node no.}$

MLP intuition

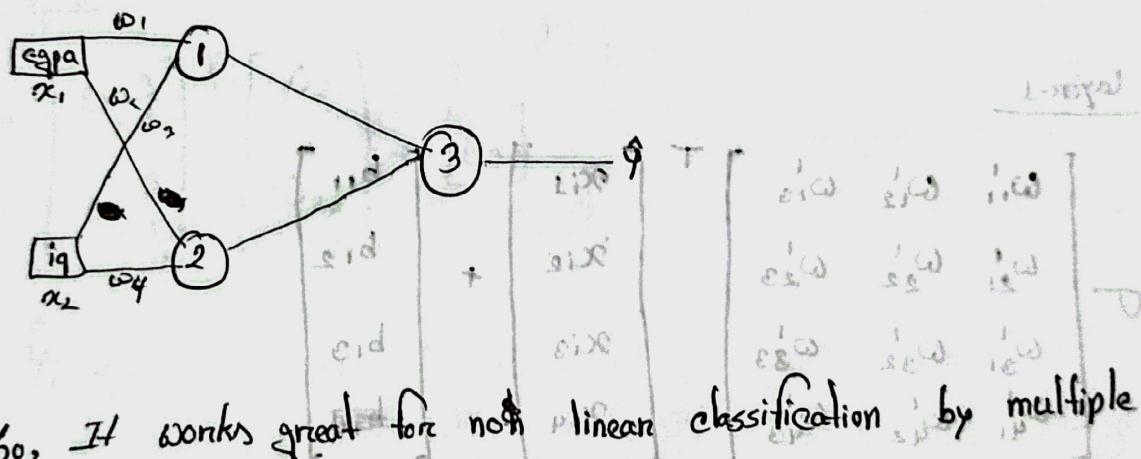


non-separable dataset



So - not linearly separable

$$(d + x^T \omega) - b = \text{constant}$$

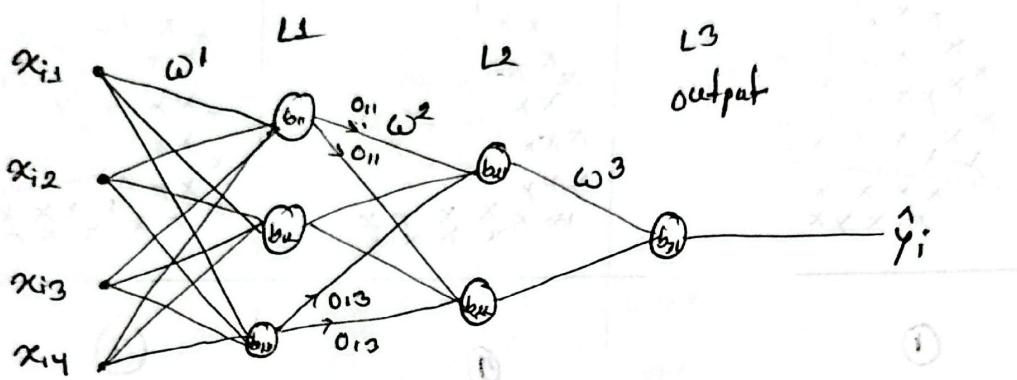


So, It works great for non linear classification by multiple linear combination perceptrons. By increasing hidden layer and adjust loss function and activation function, complex non linear dataset can be classified.

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

Forward Propagation

refined 9M



Total parameters - 26

$$\text{Perceptron} = \sigma(\omega^T x + b)$$

layer-1

$$= \sigma \begin{bmatrix} \omega_{11}^1 & \omega_{12}^1 & \omega_{13}^1 \\ \omega_{21}^1 & \omega_{22}^1 & \omega_{23}^1 \\ \omega_{31}^1 & \omega_{32}^1 & \omega_{33}^1 \\ \omega_{41}^1 & \omega_{42}^1 & \omega_{43}^1 \end{bmatrix}^T \begin{bmatrix} x_{i1} \\ x_{i2} \\ x_{i3} \\ x_{i4} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$$(4 \times 3)^T \rightarrow (3 \times 4) \quad (4,1) \cdot (3,4) + b$$

$$= \begin{bmatrix} 0_{11} \\ 0_{12} \\ 0_{13} \end{bmatrix}$$

Layer 2

$$\vec{\omega} = \begin{bmatrix} \omega_{11}^2 & \omega_{12}^2 \\ \omega_{21}^2 & \omega_{22}^2 \\ \omega_{31}^2 & \omega_{32}^2 \end{bmatrix}^T + \begin{bmatrix} 0_{11} \\ 0_{12} \\ 0_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$$= \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} \quad \left. \begin{array}{l} \text{but no factors} = N \\ \text{but no factors} = 3 \end{array} \right\} \quad \neg(\hat{x} - x) = 3$$

Layer 3

$$\alpha \begin{bmatrix} w_{11} \\ w_{12} \\ w_{21} \end{bmatrix}^T \begin{bmatrix} 0_{21} \\ 0_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \\ b_{32} \end{bmatrix}$$

+ b_{31}

nothin?

Loss Functions

Loss Function is a method of evaluating how well your algorithm is modelling your dataset.

① Mean squared error (MSE) or L2 loss

$$L = (y_i - \hat{y}_i)^2$$

$$\left| \begin{array}{l} y_i = \text{actual output} \\ \hat{y}_i = \text{predicted output} \end{array} \right.$$

~~last activation
linear~~

$$\text{cost function, } = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

① Advantages

- ⇒ Easy to interpret
- ⇒ Differentiable
- ⇒ ± local minima

Disadvantages

- ⇒ Error unit (squared)
- ⇒ Not robust to outliers.

if no outliers

② Mean Absolute error (MAE) → L1 loss

$$L = |y_i - \hat{y}_i|$$

$$c = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

if outlier

⇒ Robust to outlier

3. Huber loss

$$L = \begin{cases} \frac{\epsilon^2}{2} (\gamma - \hat{\gamma})^2 & \text{for } |\gamma - \hat{\gamma}| \leq \delta \\ \delta |\gamma - \hat{\gamma}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

hyperparameter

If dataset have outliers, it will act like MAE, otherwise MSE.

$$(\hat{\gamma})_{\text{ReLU}}(x) = (\hat{\gamma})_{\text{ReLU}}(x) - (\hat{\gamma})_{\text{ReLU}}(x) + 1$$

④ Binary Cross Entropy or log loss

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

$$C = -\frac{1}{n} \sum_{i=1}^n \gamma_i \log(\hat{\gamma}_i) - (1-\gamma_i) \log(1-\hat{\gamma}_i)$$

Adv-

* Differentiable

disadv.

* multiple local minima

* not so intuitive

$$(\hat{\gamma})_{\text{ReLU}}(x) \rightarrow \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases}$$

(normal) $\hat{\gamma}(x) \leftarrow \text{softmax} \leftarrow \text{linear} \leftarrow \text{sigmoid}$

(binary) $\hat{\gamma}(x) \leftarrow \text{sigmoid} \leftarrow \text{linear}$

(multiclass) $\hat{\gamma}(x) \leftarrow \text{softmax} \leftarrow \text{linear}$

⑥ Categorical Cross Entropy

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

$k = \text{no. of classes in data}$

For, multiple classification

last activation

Softmax

For, 3 class,

$$L = -y_1 \log(\hat{y}_1) - y_2 \log(\hat{y}_2) - y_3 \log(\hat{y}_3)$$

For that, target feature have to be one hot encoded.

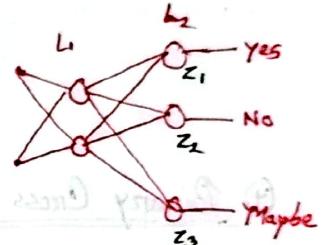
to avoid such thing,

We can use,

Sparse Categorical Cross Entropy.

It will handle internally and it is fast.

Last layer neuron number = no. of classes



$$\text{softmax} \quad f(z_i) = \frac{e^{z_i}}{\sum e^{z_i}}$$

$$f(z_i) = \frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$C = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k y_{ij} \log(\hat{y}_{ij})$$

Regression \rightarrow have outliers \rightarrow MAE [linear]
 No " \rightarrow MSE

Classification \rightarrow binary \rightarrow BCE [Sigmoid]
 Multi " \rightarrow CCE [Softmax]
 " " \rightarrow SCCE

Backpropagation

It's a method used to train neural networks, by sending errors backward through the layers to adjust the weights and bias by gradient descent to minimize loss.

Prerequisite

- Complete understanding of Gradient descent
- Forward propagation

Algorithm and steps

for i in range (epochs):

 for j in range (x.shape[0]):

 → Select 1 row from data (random)

 → Predict output using forward propagation

 → Calculate loss using loss function

 → update weights & bias using G.D.

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w}$$

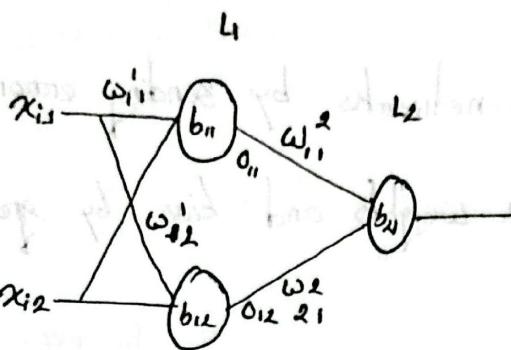
$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b}$$

Δ loss Δ loss after applying bias to loss $L(p-y) = \Delta L$

thus bias loss is loss after removing p and small diff

$$\frac{16}{16} - \frac{16}{16} = \frac{16}{16} - \frac{16}{16} = \frac{16}{16} - \frac{16}{16} = \frac{16}{16} - \frac{16}{16}$$

∴ wt and bias subject to minimization loss with one result



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

parameters = 9

Our main goal is to minimize the loss. The data inputs are constant.

If I consider \hat{y} , it is effected by weights and biases. Loss is directly effected by \hat{y} but if we want to change the \hat{y} we have to adjust weights and bias.

$$w_{\text{old}} = w_{\text{old}} - \eta \underbrace{\frac{\partial L}{\partial w_{\text{old}}}}_{\text{partial derivative}}$$

If we notice this formula, $\frac{\partial L}{\partial w_{\text{old}}}$ is the partial derivative of all the w and b . If adjust and update the values of w and b and find the best values of them.

Like, if $\text{Loss} = (y - \hat{y})^2$ and if it is the upper neural network, then, there are 9 parameters that we need to find and adjust.

$$\frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial w_{21}}, \frac{\partial L}{\partial w_{22}}, \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_{11}}, \frac{\partial L}{\partial b_{12}}, \frac{\partial L}{\partial b_{21}}$$

These are the partial contribution of weights and bias to \hat{y} .

Hence, I calculate one,

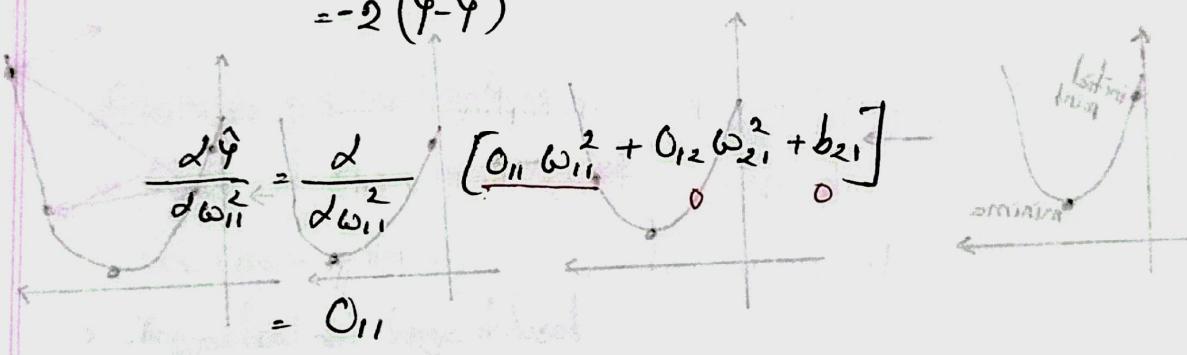
$\frac{dL}{d\omega_{11}^2}$, this is loss is not directly effect ω_{11}^2 rather, loss is effected by \hat{y} , then, \hat{y} is effected by ω_{11}^2 . It is like a chain.

That means,

$$\frac{dL}{d\omega_{11}^2} = \frac{dL}{d\hat{y}} \times \frac{d\hat{y}}{d\omega_{11}^2}$$

Now, $\frac{dL}{d\hat{y}} = \frac{d}{d\hat{y}} (y - \hat{y})^2$

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



So, $\frac{dL}{d\omega_{11}^2} = -2(y - \hat{y}) \cdot 0_{11}$,

after that, $\omega_{11}^2(\text{new}) = \omega_{11}^2(\text{old}) - \eta (-2(y - \hat{y}) \cdot 0_{11})$, new value of ω_{11}^2

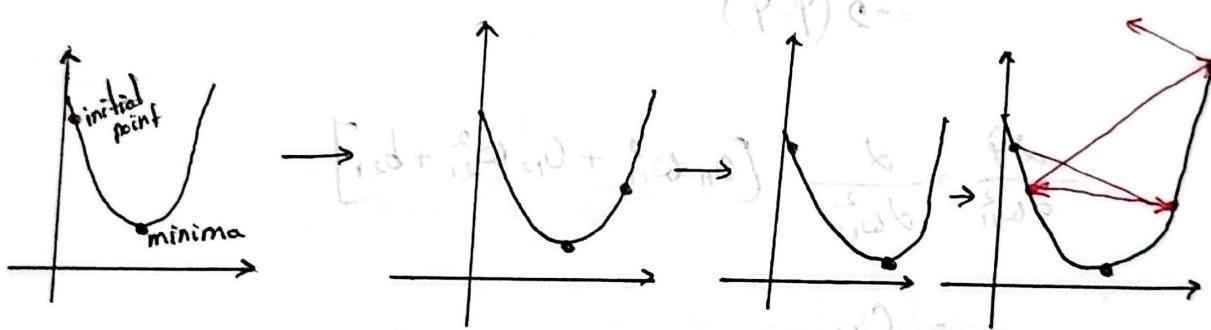
Similarly, all the other θ parameter is calculated, and it happens multiple iteration to reach the global minima.

Effect of learning rate

We introduce learning rate (η) on this formula just to get to the global minima gradually.

$$\text{If, } w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

learning rate is high, there is a possibility that, we may skip the minima. We can observe a Z pattern in gradient descent.



If, the learning rate is too less,

We may reach the minima, but it will take a lot of iteration. Means more computation and more time.

If, iteration is less, we will not reach the minima

Gradient Descent

GD is an optimization algorithm. In deep learning it is used to optimize millions of parameters (weights and bias) by minimizing a loss function.

Types

① Batch Gradient Descent

- Uses the entire data set to compute the gradient before performing a weight update.

Advantages and disadvantages

- Smooth and stable convergence
- Slow + memory heavy.
- Impractical for large dataset

② Stochastic GD (SGD)

- Update parameter after each training example.

Advantages

- Can escape local minima
- Updates are noisy
- fast update

③ Mini batch GD

- Splits dataset into small batches (e.g. 32, 64, 128) and update parameters per batch.
- Efficient on GPU
- More stable than SGD and faster than BGD
- Choosing batch size effects training efficiency.
 - batch size too small → gradient descent
 - batch size too large → stochastic gradient descent

Vanishing and Exploding Gradient

VG Gradient become extremely small as they propagate backward through layers.

Early layers learn very slowly or stop learning as gradient small

Why

\Rightarrow Mostly happens if its a very deep network.

\Rightarrow repeated multiplication of derivative < 1 (Sigmoid/tanh)

\Rightarrow slow or stop convergence.

Solution

\Rightarrow Use ReLU or variants

\Rightarrow Use Batch normalization

\Rightarrow Carefull weight initialization

EG

Oposite of vanishing Gradient.

Why

\Rightarrow High learning rate

\Rightarrow Improper weight

Solution

\Rightarrow Gradient clipping

\Rightarrow Lower learning rate

\Rightarrow Normalization layer

Fine Tuning Hyper-parameter

Layers and Neurons

more layers → capture complex patterns

more neurons → increase representation capacity (more than sufficient)

Batch size

→ Small batch size (e.g. 32, 64, 128) → preferable better perform and generalization

→ Large batch → faster but not so generalize. May cause overfitting.

Number of epoch

→ Use early stopping to avoid overfitting

Early Stopping

Early Stopping \Rightarrow a regularization technique used to prevent overfitting. Training stops before model starts overfitting.

\Rightarrow Used to stop epoch before overfitting.

Code

callback = Early Stopping (

monitor = "val-loss"

`model.fit(x-train, y-train, validation-data=(x-test, y-test), epoch=3000, callbacks=[callback])`

bitsaw at last year split out of (q-e-b) start hognub. The

2000 wooden combs)

800-02-000A

CHM-302-005

100-05 ← 149

La-ca-ca-ca-ca-ca la-ca-ca-ca-ca-ca la-ca-ca-ca-ca-ca

1) Introducing a line

(raib.tupi, "uler" - mitariloo, 25) was

• 15 •

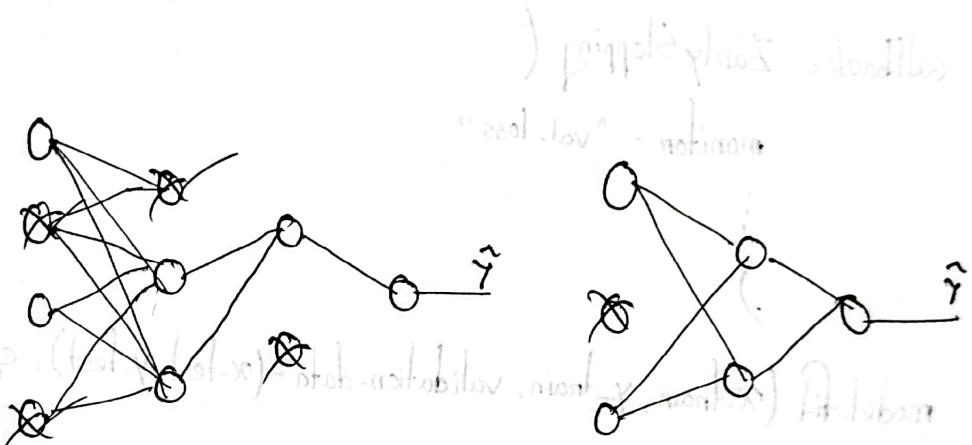
卷之三

* (2-3) beyond

Dropout

Dropout is a regularization method used to reduce overfitting in neural network. It works by randomly dropping or deactivate a subset of neurons during training forcing the network to learn redundant, robust representation.

→ force each neurons to learn useful, general features independently.



If, dropout rate ($d = 1-p$) is too high, may tend to underfit.

Common dropout rate

ANN \rightarrow 10 - 50%

CNN \rightarrow 40 - 50%

RNN \rightarrow 20 - 30%

from tensorflow.keras.layers import Dropout

```
model = Sequential([
    Dense(128, activation = "relu", input_dim=),
    Dropout(0.2),
    Dense(),
    Dropout(0.2),
])
```

Regularization

It's a set of techniques used to reduce overfitting and improve generalization of machine learning and deep learning models. It introduces constraints or penalties to control model complexity, encouraging model to learn simple and generalize patterns. It reduces the values of weights and bring them closer to zero.

$$C = L(\hat{y}, y) + \underbrace{\frac{\lambda}{2n} \sum_{i=1}^k \|w_i\|^2}_{\text{penalty}} \rightarrow L_2$$

For, L_1 ,

$$\frac{\lambda}{2n} \sum_{i=1}^k \|w_i\|$$

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

For, a sample, loss function will be,

$$L' = L + \frac{\lambda}{2} \sum \|w_i\|^2$$

Now, $\frac{\partial L'}{\partial w_0} = \frac{\partial L}{\partial w_0} + \lambda w_0$

$$\therefore w_n = w_0 - \eta \left(\frac{\partial L}{\partial w_0} + \lambda w_0 \right)$$

$$= (1 - \eta \lambda) w_0 - \eta \frac{\partial L}{\partial w_0}$$

$w \downarrow$

model = Sequential()

model.add(Dense(128, input_dim=2, activation="relu", kernel_regularizer=tensorflow.keras.regularizers.l2(0.005)))

L1

Activation Function

It defines how the weighted and biased sum of input + bias is transformed before passing to the next layer in a neural network.

$$Z = \omega^T x + b$$

activation function, $\alpha = f(z)$

Ideal activation Function

- 1) Non linear
 - 2) Differentiable
 - 3) Computationally inexpensive
 - 4) Zero centered ($\text{mean} = 0$) (tanh)
 - 5) Non-saturating (relu)

Sigmoid

$$f(x) = \frac{1}{1+e^{-x}}$$

$$f'(x) = f(x)(1-f(x))$$

range : (0, 1)

- Commonly used in output layer

(U1a) few words before

Advantages

- ⇒ Smooth and differentiable
- ⇒ Useful for probabilistic interpretation
(Binary classification)
- ⇒ Non-linear

Disadv

- ⇒ Vanishing gradient due to saturation
- ⇒ Not zero centered
- ⇒ Slow convergence

Hyperbolic Tangent (tanh)

$$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = \frac{1-f(x)}{1+f(x)}$$

Range : (-1, 1)

- Hidden layers in early RNN and MLPs.

Adv

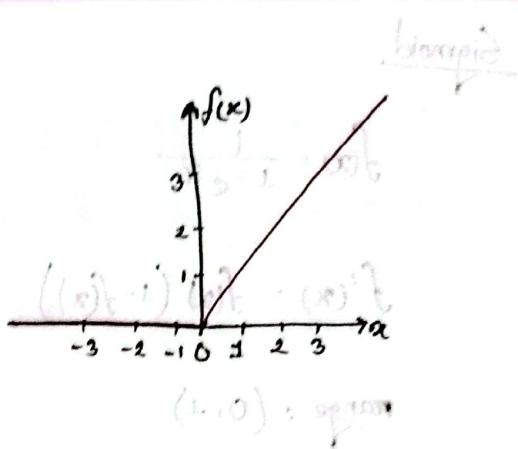
Zero centered others like sigmoid

disadv

- ⇒ Vanishing gradient

Rectified Linear Unit (ReLU)

$$f(x) = \max(0, x)$$



$$f(x) = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

Advantages

- Computational efficiency
- Non linear
- Mitigates vanishing gradient
- Introduces sparsity (Some neurons deactivate)

Disadv

Dying ReLU problem: neurons can inactive (output 0), if weight update incorrectly.

Non zero centered

Common

Hidden layers in ANN, RNN, CNN, MLP.

Dying ReLU

It occurs when a neuron's output becomes 0 for all input, and never recover during training. Means,

⇒ Neurons stop learning

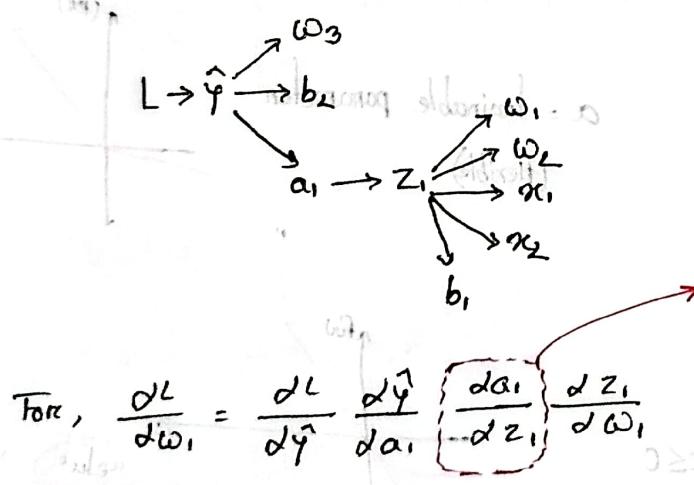
⇒ Becomes inactive

Why happens

During training, weights update.

$$\omega_n = \omega_0 - \eta \frac{\partial L}{\partial \omega_0}$$

Dependency



If, this derivation became zero, then there will be no weight update. That's when we will face dying ReLU.

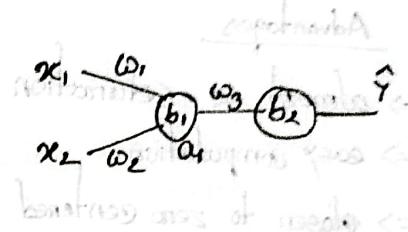
When happens

⇒ Large learning rate

⇒ High neg. bias

ReLU

$$(s, s, 0, 0) \times \text{ReLU} = (0)^4$$



ReLU

$$a_i = \max(0, z_i)$$

$$z_i = \omega_1 x_1 + \omega_2 x_2 + b_i < 0$$

$$0 < x_i$$

$$0 > x_i$$

$$\therefore a_i = 0$$

$$30, \frac{\partial a_i}{\partial z_i} = 0$$

(tiny noise taken into account)

$$0 < x_i \quad x_i \quad 0 > x_i \quad (1-x_i) > 0$$

(big noise taken into account)

$$0 < x_i \quad x_i \quad 0 > x_i \quad 0 > x_i$$

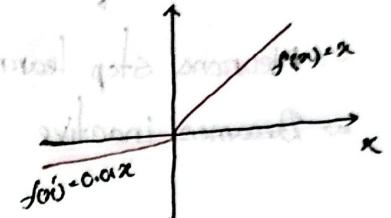
Variation of Relu

Most of them solve dying relu problem.

Leaky Relu

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

Typically,
 $\alpha = 0.01$
(reached) value



activation = LeakyRelu(alpha=0.01)

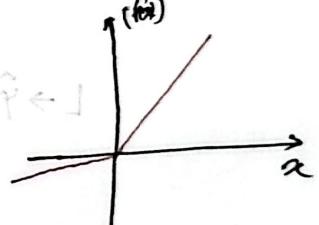
Advantages

- ⇒ almost no saturation
- ⇒ easy computation
- ⇒ closer to zero centered

Parametric Relu

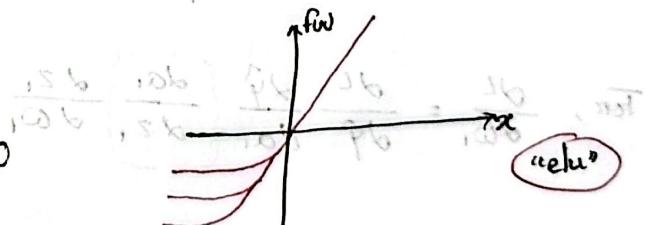
$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x \leq 0 \end{cases}$$

α = trainable parameter
(flexible)



Elu (Exponential Linear unit)

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



Selu (Scaled Elu)

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

$$\alpha \approx 1.673$$

$$\gamma = 1.051$$

shortest gradient update time
stable for higher α

Weight initialization

What not to do, estimate change primarily with respect to gradient flow

- ① All weights are zero
- ② Are constants
- ③ Random small (-0.01 to 0.01)
- ④ Random large. (-3 to 3)

Initializations to maintain stable, stochastic optimization so gradient flow

What can be done

Xavier/Glorot initialize

number does not

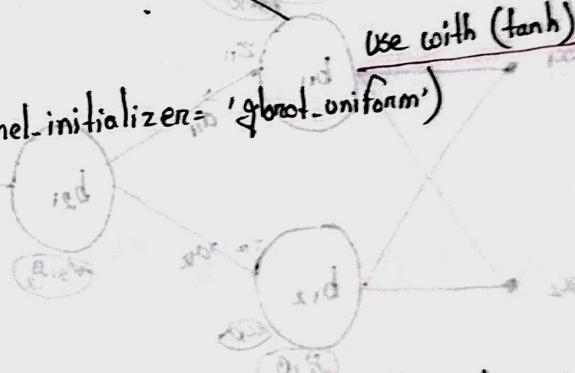
transformation

$$\text{np.random.randn(dim)} * \frac{1}{\sqrt{\text{fan-in}}}$$

fan-in = number of inputs coming to the node

use with (tanh), sigmoid

Dense(256, activation = 'tanh', kernel_initializer = 'glorot-uniform')



He initialize

$$\text{np.random.randn(dim)} * \sqrt{\frac{2}{\text{fan-in}}} \leftarrow (N(0, S)\beta \leftarrow N(0, S) \leftarrow N(0, S) \leftarrow S$$

use with relu and variant

(selu, leakyrelu)

model = Sequential([

Dense(256, activation = 'relu', kernel_initializer = 'he-normal')])

initialization config

start forward pass -- what is

using update add the

gradient regularization layers to help

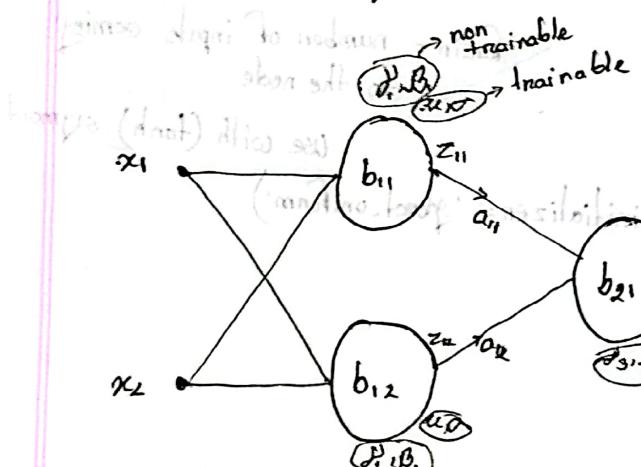
((0 - minibatch_size * gradient) / size) * learning_rate

((1 / minibatch_size) * size) * learning_rate

Batch Normalization

It is a technique to improve the training speed, stability, and performance of deep learning networks. It normalize the input of each layer so that they have a mean of 0 and variance of 1, within a mini-batch.

During training, as parameter update, the distribution of activations in each layer changes - it is called Internal Covariate Shift.



For, each neuron,
 $w, b \rightarrow$ different
 $\mu, \sigma \rightarrow$ different

$$z_{11} = w_{11}x_1 + w_{21}x_2 + b_{11}$$

rest will be same process

Process
 $z_n \rightarrow z_{nN} \rightarrow z_{nBN} \rightarrow g(z_{nBN}) \rightarrow a_{nBN}$

$$z_{nN} = \frac{z_n - \mu_B}{\sigma_B + \epsilon} \quad | m = \text{batch size}$$

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

variance, $\sigma_B = \sqrt{\frac{1}{m} \sum_{i=1}^m (z_i - \mu_B)^2}$

$$z_{nBN} = \gamma z_{nN} + \beta$$

gives flexibility

Adv.

- ⇒ Stable — Wide range hyperparameters
- ⇒ Faster — higher learning rate
- ⇒ act like regularizer
- ⇒ Impact of weight initialization reduce

Code

```
- model.add(Dense(3, activation='relu', input_dim=2))
model.add(BatchNormalization())
```

Optimizers

(MATH) for neural networks optimization

These algorithms used to update parameters of a neural network in order to minimize the loss function. It helps more giving good outcome.

Till now, we use Gradient descent and its types.

There were many problem with these optimizers like,

1> Choosing right learning rate

2> get stuck in local minima

3> Saddle point where gradient $\frac{\partial L}{\partial \omega}$ is 0 on descent $\leftarrow (0,0)$ is saddle point $\leftarrow (0,0-1,0)$ need not be

Optimizers that we will learn

1> Momentum

2> AdaGrad

3> RMSprop

4> NAG

5> adam



$$(\text{Adam} \cdot (\text{SGD} + \text{Adagrad}))_{\text{new}} = (\text{Adam} \cdot \text{SGD})$$

Exponentially weighted moving avg (EWMA)

It's a technique to compute a smooth average of a sequence of numbers by giving more weight to recent observation and less to older ones.

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

V_t = current moving avg

V_{t-1} = previous " "

x_t = current value

β = decay rate or smooth factor
 $(0 < \beta < 1)$

β determines how much weight is given to the past vs the present.

\Rightarrow higher β (0.9) \rightarrow smoother, slow react to present

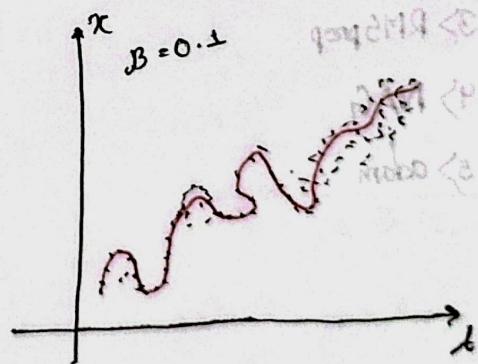
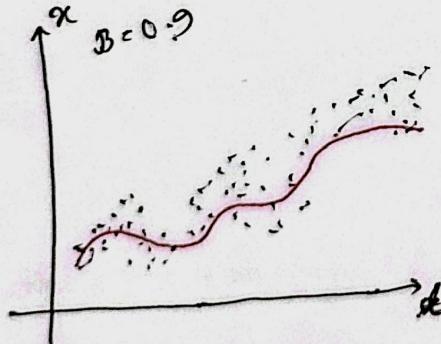
\Rightarrow for lower ($0.1-0.5$) \rightarrow opposite

D_t	α_t (Temp)
D_1	43°
D_2	33°
D_3	25.5°
D_4	30.6°

Observation

$$V_t = (1-\beta) (x_t + \beta x_{t-1} + \beta^2 x_{t-2} + \beta^3 x_{t-3} + \dots)$$

$$\beta = 1 - \alpha$$



$[x_t = df[::].ewm(alpha=0.9).mean()]$

Momentum

An optimization technique used to accelerate gradient descent by taking account the past gradient.
It behaves like a ball rolling down a hill - gaining 'momentum' as it moves.

Gradient descent,

$$\omega_{t+1} = \omega_t - \eta \nabla \omega_t$$

In momentum,

$$\omega_{t+1} = \omega_t - v_t$$

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

$\beta \rightarrow$ momentum coefficient
normally (0.9)

Adv

- 1) faster convergence
- 2) smooth updates
- 3) Escape local minima

Disadv

- 1) May overshoot if β and η is too high
- 2) Requires hyper tuning

Nesterov Accelerated Gradient (NAG)

An improvement of Momentum. It adds a look ahead step - instead of computing the gradient at the current position, it looks ahead in the direction of the accumulated momentum, before computing the gradient.

$$\nabla f(\theta) = \text{backward}$$

backward gradient \rightarrow $\nabla f(\theta)$

$$\nabla f(\theta) = \text{backward}$$

$$\nabla f(\theta) + \lambda \nabla \theta = \nabla f(\theta)$$

forward

look ahead gradient \rightarrow forward gradient

forward gradient \rightarrow backward gradient \rightarrow forward gradient