

Date: 26 April 2020

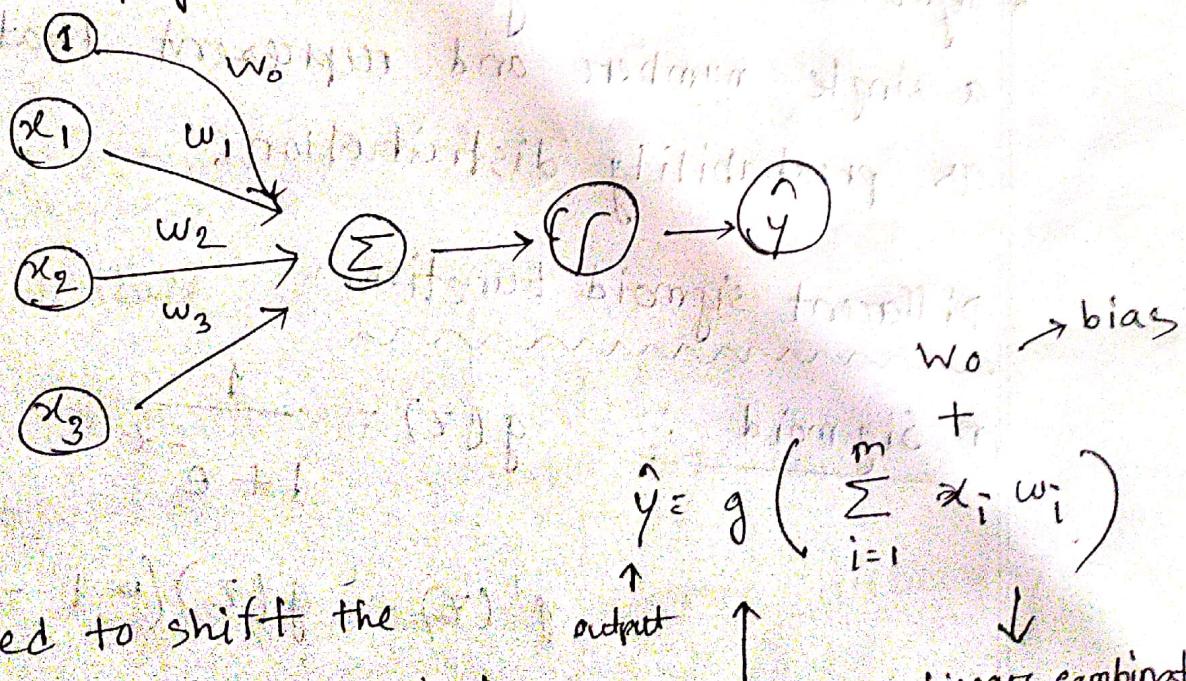
- * Always focus on documentation
- * Linear algebra (matrix multiplication)

Author: Kazi Amit Hasan
Email: kazi.amithasan.89@gmail.com

Introduction to deep learning : (MIT 6.S191)

- The Perceptron
- The structure building block of deep learning

Forward Propagation



bias: used to shift the activation function left/right based on inputs

output
non-linear activation function

Linear combination of inputs

$$\hat{y} = g(\omega_0 + \mathbf{x}^T \mathbf{w})$$

↑ activation function

where, $\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ $\mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_m \end{bmatrix}$

Ex: Sigmoid function,

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

- * Sigmoids are really useful when we want a single number and represent that number as probability distribution.

Different sigmoid function:

① Sigmoid : $g(z) = \frac{1}{1 + e^{-z}}$

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

⑩ Hyperbolic tangent:

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

$$g'(z) = 1 - g(z)^2$$

⑪ Rectified Linear unit (ReLU)

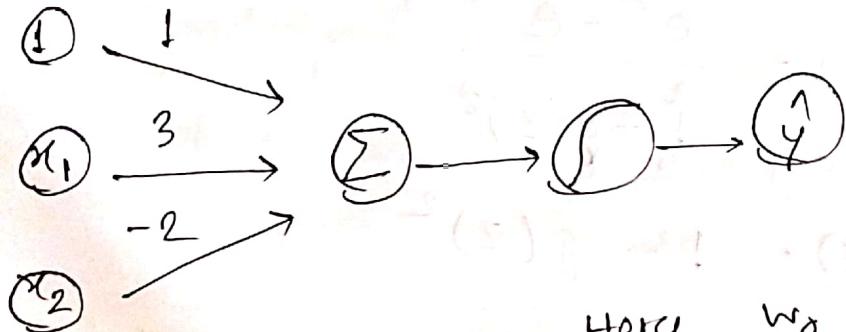
$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1 & , z > 0 \\ 0 & , \text{otherwise} \end{cases}$$

Importance of activation function:

- to introduce non-linearities into the network.
- x Linear activation function produce linear decisions
- x no matter what the network size is.
- x Non-linearities allow us to approx. arbitrary complex functions.

Examples

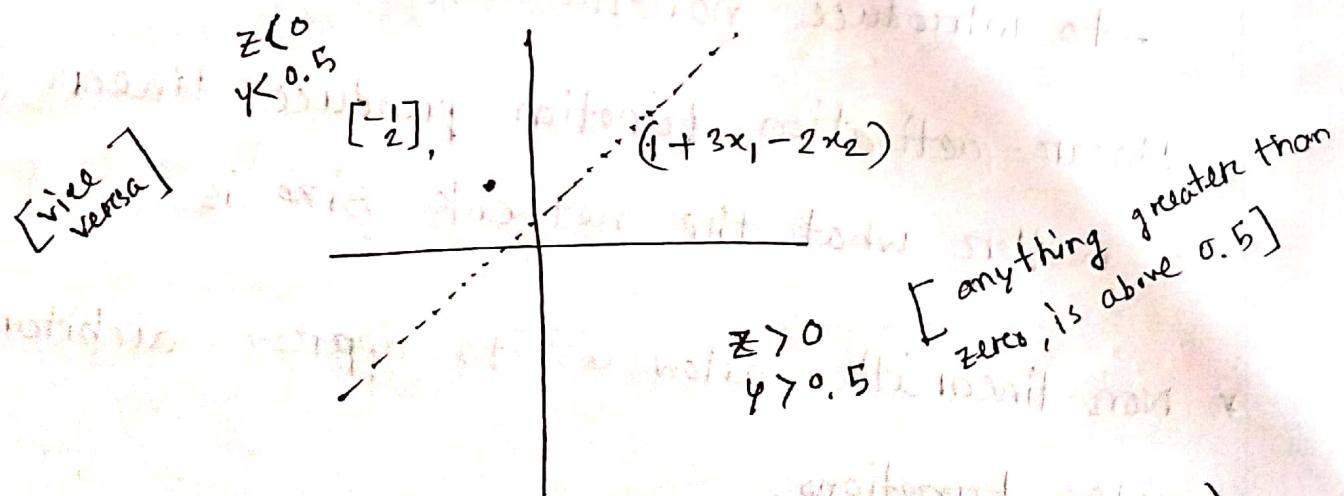


$$\text{Here, } w_0 = 1$$

$$w = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{x}^T w) \\ &= g(1 + \mathbf{x} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix})\end{aligned}$$

$$\therefore \hat{y} = g(1 + 3x_1 - 2x_2) \quad [2D \text{ line}]$$

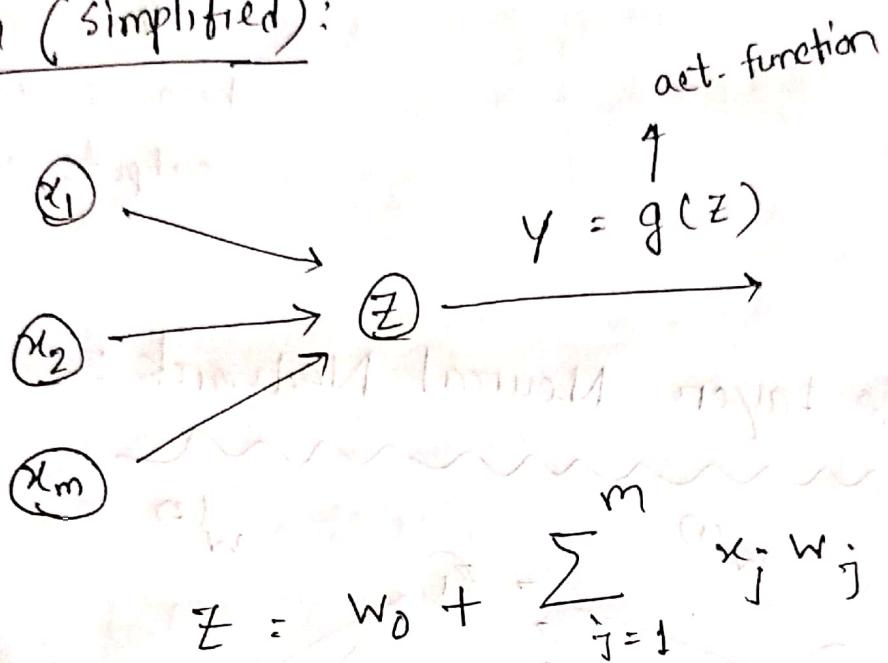


Protivias

$$\begin{aligned}\text{Assume, } \mathbf{x} &= \begin{bmatrix} -1 \\ 2 \end{bmatrix} & \therefore \hat{y} &= g(1 + 3(-1) - (2)(2)) \\ &= g(-6) \approx 0.002 & [\text{sigmoid}]\end{aligned}$$

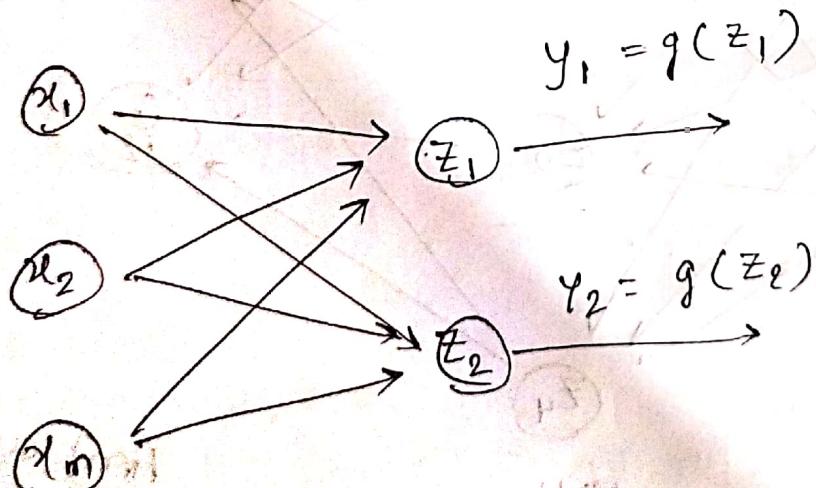
Building Neural Networks with Perceptron:

Perceptron (simplified):



Multi-output Perceptron:

$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

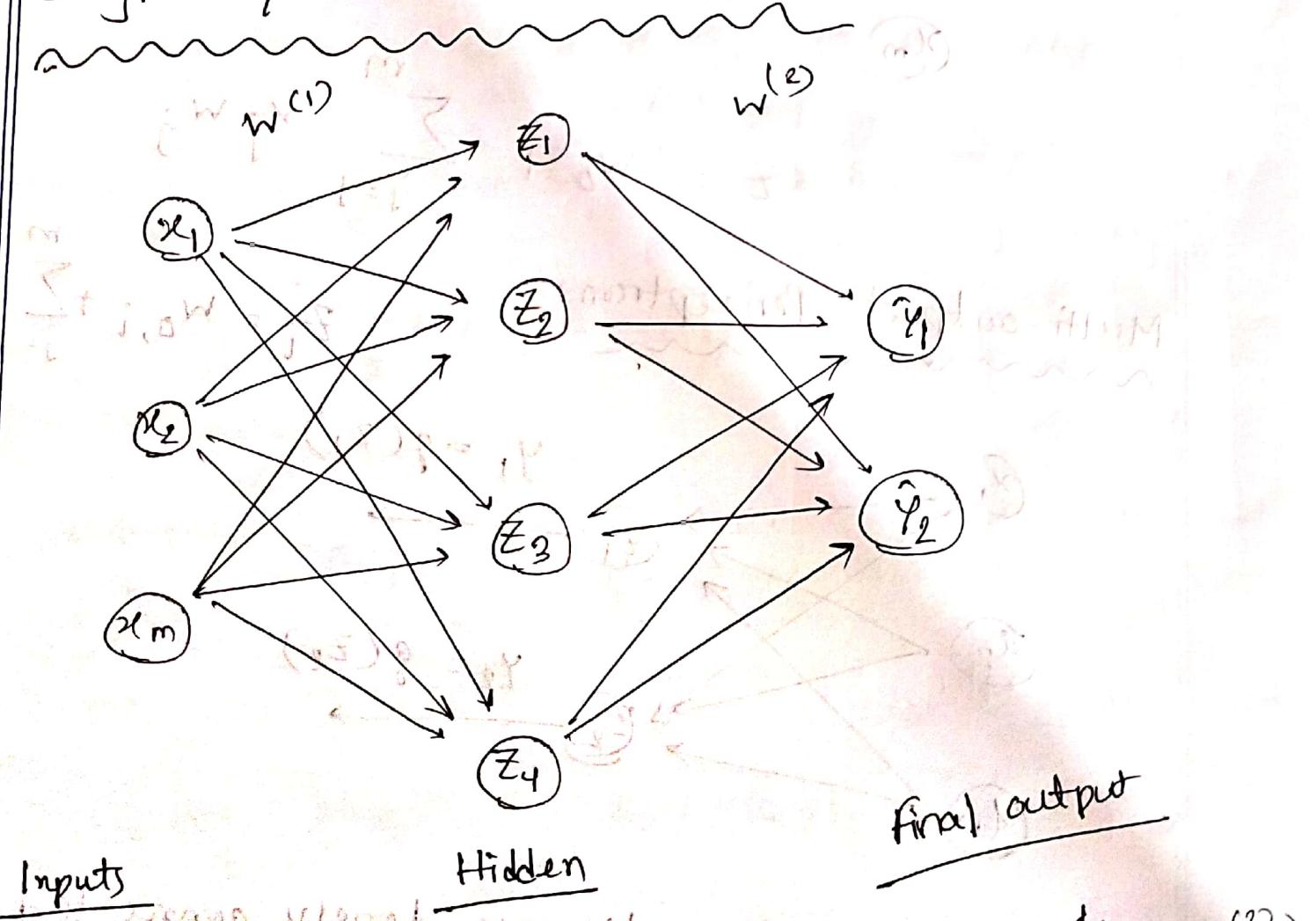


Dense layer: all inputs are densely connected to all outputs.

```
import tensorflow as tf  
layers = tf.keras.layers.Dense(units=2)
```

because we have two outputs in the layer

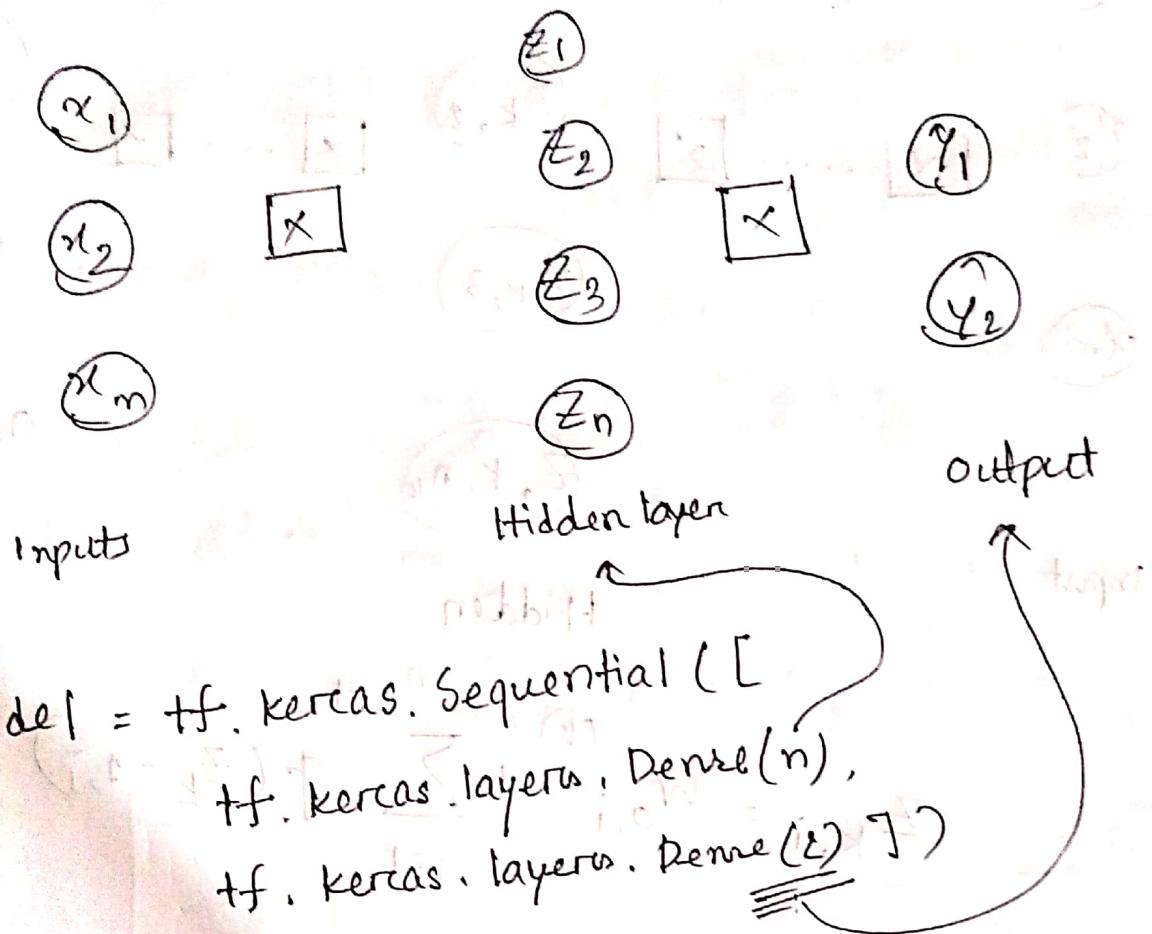
Single Layer Neural Network :



$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)}$$

$$y_i = g(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)})$$

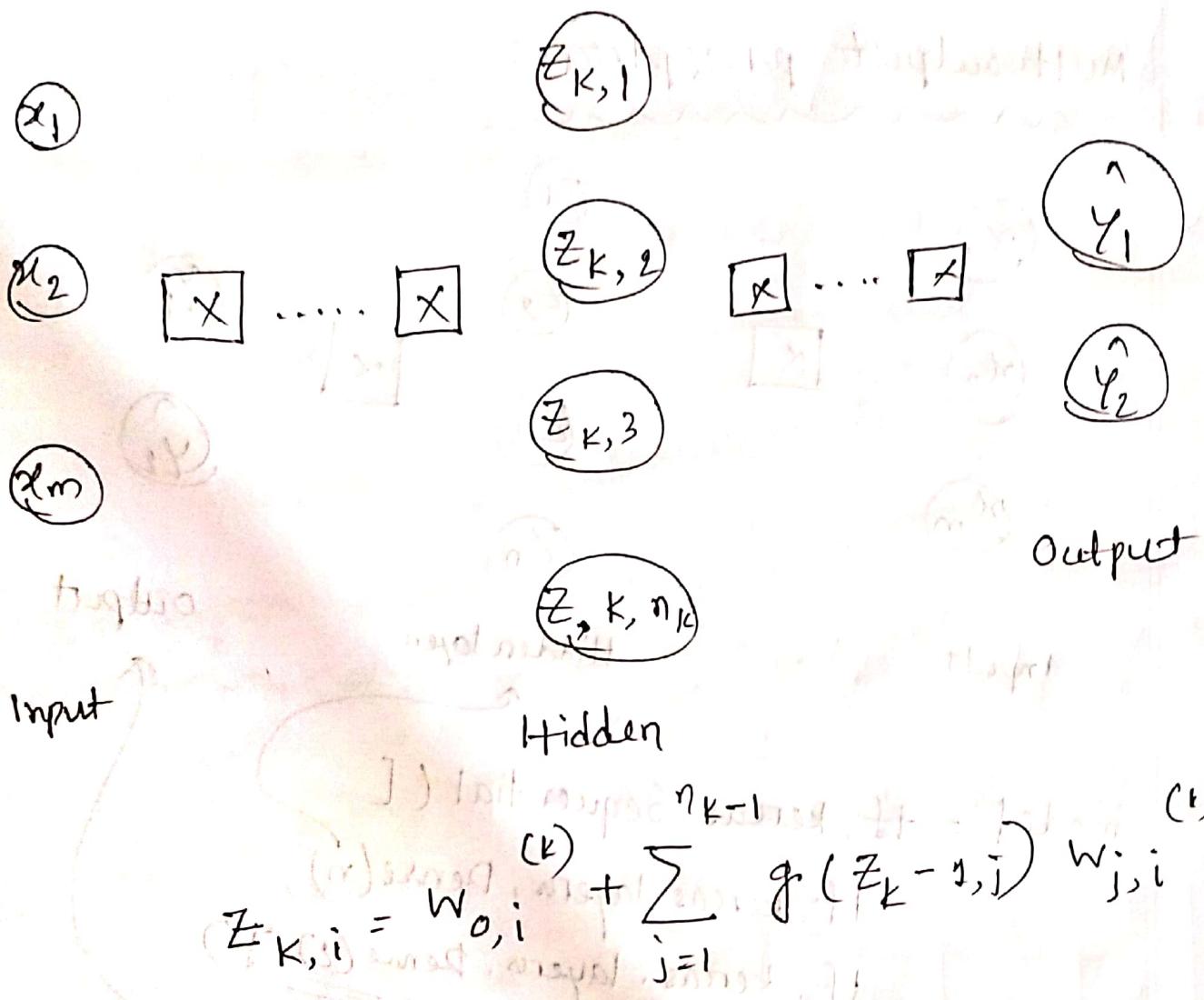
Multi output perceptron :



model = tf.keras.Sequential [
 tf.keras.layers.Dense(n),
 tf.keras.layers.Dense(2)])

Deep Neural Network:

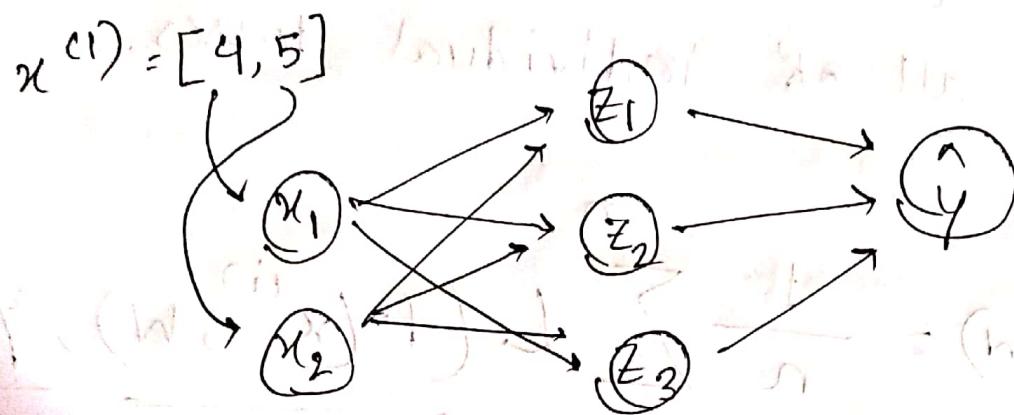
P.T.O.



tf.keras.layers.Dense(n_1),
 tf.keras.layers.Dense(n_2),
 : : : :
 tf.keras.layers.Dense(2)

Applying Neural Networks:

Ex: Will I pass? x_1 = No. of lectures attended
 x_2 = Hours spent on project



$x^{(i)} = [4, 5]$ pred = 0.1
which is wrong
Actual = 1 (Yes, passed)
train याते रथा

Quantifying Loss:

- the loss of network that measures that cost incurred from incorrect predictions.

$$L(f(x^{(i)}; w), y^{(i)})$$

pred actual

Empirical Loss: Obj function / Cost function / Empirical Risk

The empirical loss measures the total loss over entire dataset.

- Avg of all of individual losses.

$$J(w) = \frac{1}{n} \sum_{i=1}^n L(f(x^{(i)}; w), y^{(i)})$$

$f(x)$	x	y
0.1	X	1
0.8	X	0
0.6	✓	-

Binary Cross Entropy Loss

$\text{loss} = \text{tf.reduce_mean}(\text{tf.nn.softmax_cross_entropy_with_logits}(\text{logits}, \text{labels}))$

entropy-with-logits(y , predicted)

Cross entropy loss can be used with models that's output is between 0 & 1.

$$J(w) = \frac{1}{n} \sum_{i=1}^n y^{(i)} \log \left(\frac{f(x^{(i)}; w)}{\text{pred}} \right) + (1 - y^{(i)}) \log \left(\frac{1 - f(x^{(i)}; w)}{\text{pred}} \right)$$

actual

What if we have a regression problem instead to classification problem? [grade for 20]

Mean Squared Error loss:

used to regression model that output continuous real numbers.

Prediction of network, actual true output, then subtracting them, take the squared error and that's the loss, this is the network should optimize.

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y^{(i)} - \underline{f(x^{(i)}; w)})^2$$

actual pred

Training Neural Networks:

Loss Optimization:

- we need to find the network that achieves the lowest loss.

$$\underline{w} = \arg \min \frac{1}{n} \sum_{i=1}^n l(f(x^{(i)}; w), y^{(i)})$$

$$\underline{w} = \arg \min J(w)$$

set of weights that tries to minimize the loss function.

- * Loss function is a function of the network weights

compute gradient = $\frac{\partial J(w)}{\partial w}$

↳ we are taking this to understand the direction of maximum ascent, basically tells us the conditions of the path (up/down) → up direction $\frac{\partial J(w)}{\partial w} > 0$, reverse if $\frac{\partial J(w)}{\partial w} < 0$, until it converge to a lowest point - local minima.

This is gradient descent algorithm

Algo: ① Initialize weights randomly

② Loop until convergence

③ Compute gradient = $\frac{\partial J(w)}{\partial w}$

④ Updates weights, $w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$

⑤ Return weights

learning rate

while True:

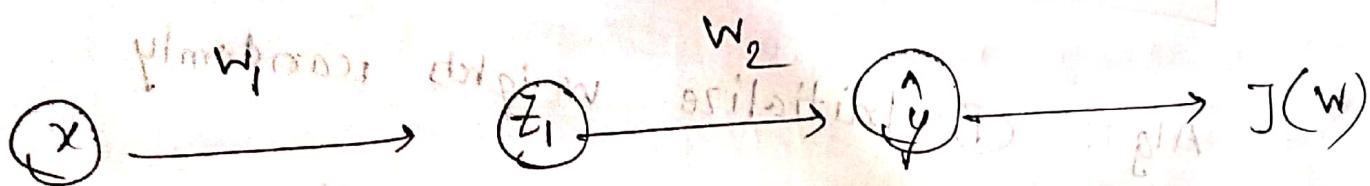
 with tf.GradientTape() as g:

 loss = compute_loss(weights)

 gradient = g.gradient(loss, weights)

 weights = weights - lr * gradient

→ back-propagation : (application of chain rule)



How a small change in one weight (w_2/w_1)

can affect the final loss?

$$\frac{\partial J(w)}{\partial w_2} = \frac{\partial J(w)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_2} \quad [\text{chain rule}]$$

Again,

$$\frac{\partial J(w)}{\partial w_1} = \frac{\partial J(w)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial w_1}$$

$$= \frac{\partial J(w)}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_1} \times \frac{\partial z_1}{\partial w_1}$$

[chain rule]

→ Repeat this for every weight in the layer
using gradients from latter layers.

Neural Networks in Practice!

Optimization:

Loss functions can be difficult to optimize.

optimizing through gradient descent:

$$w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$$

learning rate

(determines how much steps we should take in the direction of the gradient)

small learning rate converges slowly

and get stuck in false local minima

- we may stuck in local minima

because we are not as aggressive we should be

- Large learning rate overshoot, become unstable and diverge
- Smooth learning rate converges smoothly and avoid local minima.

How to deal with it?

- ① Try a lot of different learning rates.
- ② Design adaptive learning rate that adapts to the landscape.

Adaptive learning rates:

- means learning rates are no longer fixed.
- it can increase/decrease through training.
- depending on:
 - ① how large the gradient is
 - ② how fast the learning happening
 - ③ size of particular weights

Algorithms:

- ① SGD
- ② Adam
- ③ Adadelta
- ④ Adagrad
- ⑤ RMSProp

optimizer = tf.keras.optimizer.SGD()

Mini-batch:

- Backpropagation is very expensive to compute

Now we taking a single point instead of all and compute the gradient, then in next iteration, we take a different point.

S to chastic Gradient Descent:

Algorithm:

1. Initialize weights w randomly

2. Loop until converge

3. Pick a single point i

4. Compute gradient, $\frac{\partial J_i(w)}{\partial w} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(w)}{\partial w}$

5. Update weights, $w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$

6. Return weights.

Downsides :

① Single point \rightarrow very noisy

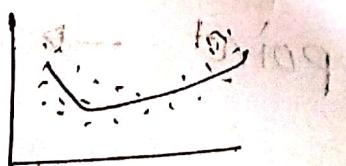
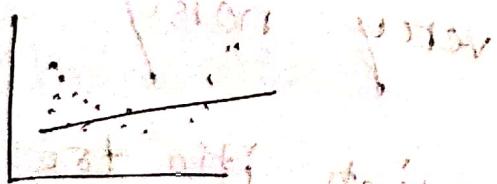
② If we consider all points, it's too expensive.

Solution:

Take a Minibatch

- more accurate estimations of gradient
- smoother convergence
- allow for ~~real~~ larger learning rates
- Mini-batches leads to fast training
- can parallelize computations + achieve significant speed increase on GPU's

Overshooting: (generalization in ML)



Underfitting

Ideal fit

Overfit

- model does not have capacity to fully learn data

Too complex.

Regularization:

- it is a technique that constrain our optimization problem to discourage complex models.
- it improves generalization of our model on unseen data.

→ Dropout:

- during training, some of the activations are set to zero. (randomly).
tf. keras.layers.Dropout($p=0.5$)
- forced network to not rely on any 1 node.

Date : _____ / _____ / _____

Early stopping:

- stop training before we have a chance to overfit.

Thanks