

# TENSOR FLOW BASIC PRINCIPLE

- \* Define a graph of Computations to perform and then TensorFlow takes that graph and runs it efficiently using optimized C++ Code
- \* It is possible to break up the graph into several chunks and run them in parallel across multiple CPU or GPUs
- \* TensorFlow developed by GOOGLE BRAIN TEAM
- \* TensorFlow "autodiff" feature can automatically and efficiently compute the gradients for you
- \* TensorFlow operations (OPS) can take any number of Input and produce any number of outputs. For example, the addition and multiplication ops each take two inputs & produce one output.

Eg:-

```
training-op = tf.assign(theta, theta - learning-rate * gradients)
```

- \* Constants and Variables take no input (source ops), the Input and outputs are multi dimensional array called tensors.

Feeding Data to the Training Algorithm:

- \* placeholder These nodes are special because they don't actually perform any computation, they just output the data you tell them to output at runtime.

4D-Tensor

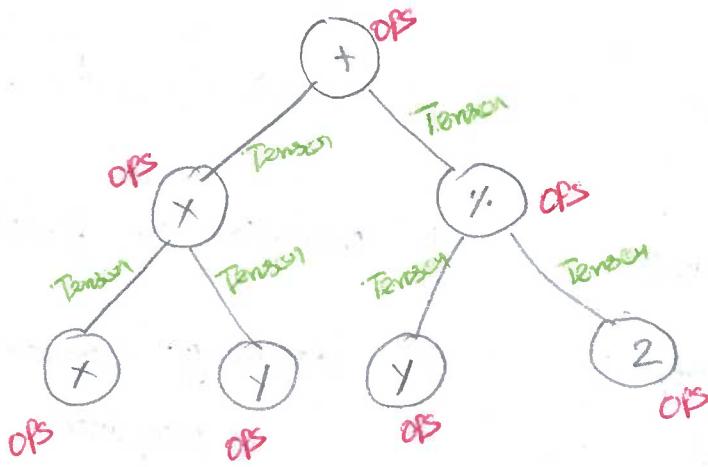
- \* [mini-batch size, Height, Width, Channels]

140      256      256      1

How to create a Tensor?

- 1) operations (ops)
- 2) Tensor

The operation is the one that results in the tensor

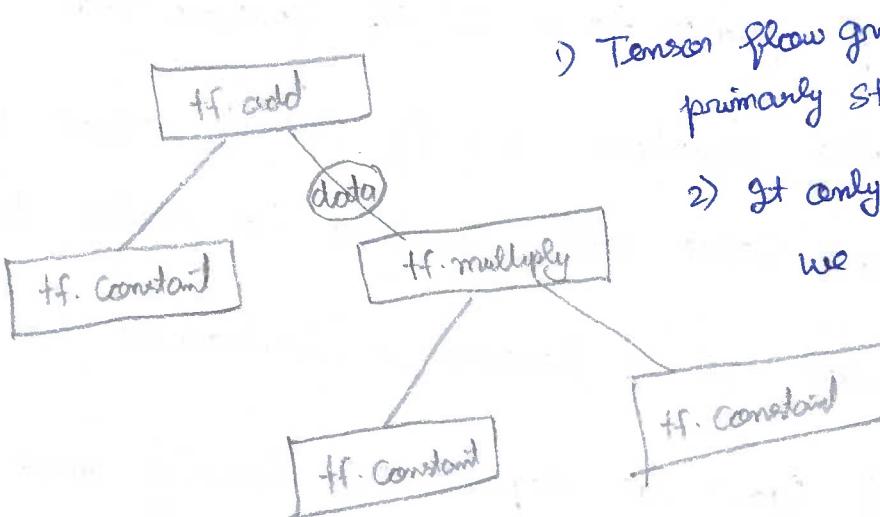


In the Computational graph

- 1) OPS are the nodes
- 2) Tensors are the edges

Hello = tf.constant(name='op1', value='Hello world')

Tensor                          OPS



- 1) Tensor flow graph doesn't primarily store data
- 2) It only yield result when we execute `sess.run` along with the input arguments

## Lazy Execution

- 1) execution on Session()
  - 2) explicit session scope is required

## Eager Execution

- 1) execution on function call
  - 2) defining Session is no longer required

0-tensor (SCALAR)

1- tensor (VECTOR)

2-tensor (matrix)

## Random Tensors.

Many machine learning algorithms learn by

performing updates to a set of tensors that hold weights.

These update equations usually satisfy that weights initialized at the same value will continue to evolve together. Thus, if the initial set of tensors is initialized to a constant value the model won't be capable of learning much.

tf. expand\_dims → Convert a vector into a row vector or column vector

tf. squeeze → Remove all dimension of size 1 from a tensor  
(Convert row or column vector into a flat vector)

## TENSOR FLOW GRAPHS:

Any Computation in Tensorflow is represented as an instance of `tf.Graph` object. Such graph consist of a set of instances of `tf.Tensor` objects and `tf.Operation` objects

### What is tf.Operation?

A call to an operation like "tf.matmul" creates a "tf.Operation" instance to mark the need to perform the matrix multiplication operation.

## TENSOR FLOW VARIABLE

Stateful Computation → Learning algorithms are essentially rules for updating stored tensors.

`tf.Variable()` class provides a wrapper around tensors which allows for stateful computation

[NOTE] Variables have to be explicitly initialized

How to update the value of an existing Variable?

`tf.assign()`

## GRAPH, SESSION and FETCHES

- \* Tensor Flow Involves two main Phases
  - 1) CONSTRUCTING A GRAPH
  - 2) EXECUTING THE GRAPH
- \* A **Session** object is the part of the Tensor-Flow API that communicates between Python objects and ~~at~~ data on our end, and the actual Computational System where memory is allocated for the objects we define, intermediate variables are stored and finally results are fetches for us.

## OPTIMIZER:

The next thing we need to figure out is how to **MINIMIZE** the **Loss function**

optimizers update the set of weights Iteratively in a way that decreases the loss over time.

The Most Commonly Used approach is GRADIENT DESCENT

- \* While Convergence to the GLOBAL MINIMUM is guaranteed for CONVEX FUNCTION, for NON-CONVEX problems they can get Stuck in the LOCAL MINIMA

Difference between CONVEX and NON-CONVEX OPTIMIZATION?

A Convex optimization Problem  
maintains the property of linear  
programming problem

A NON-CONVEX problem maintains the property of non-linear programming problem

There can be only one optimal solution, which is globally optimal  
as you might prove that there is no feasible solution to the problem

They can have multiple locally optimal points and it can take a lot of time to identify whether the problem has no solution or if the solution is global.

It might lead you to  
DEAD END

`tf.name_scope()`: This is used to group together parts that are related to inferring the output, like defining the loss, and setting and creating the training object.

Input Size: [None, 28, 28, 1]

- 1) None: we have an unknown number of Images
- 2) each  $28 \times 28$  pixel
- 3) 1 Color channel (since these are gray scale Images)

### Two Popular Example of GAN IMPLEMENTATION

- 1) Deep Convolutional GAN (DCGAN)
- 2) Conditional GAN (CGAN)

GAN can generate **NEW** Celebrity faces that are not of real people by performing latent space interpolations

GANs are able to learn how to model the Input distribution by training two competing (and cooperating) networks referred to as **GENERATOR** and **DISCRIMINATOR** (sometimes known as CRITIC).

**GENERATOR:** the role of the generator is to keep figuring out how to generate fake data or signals that can fool the discriminator.

**DISCRIMINATOR:** Trained to distinguish between fake and real signals.

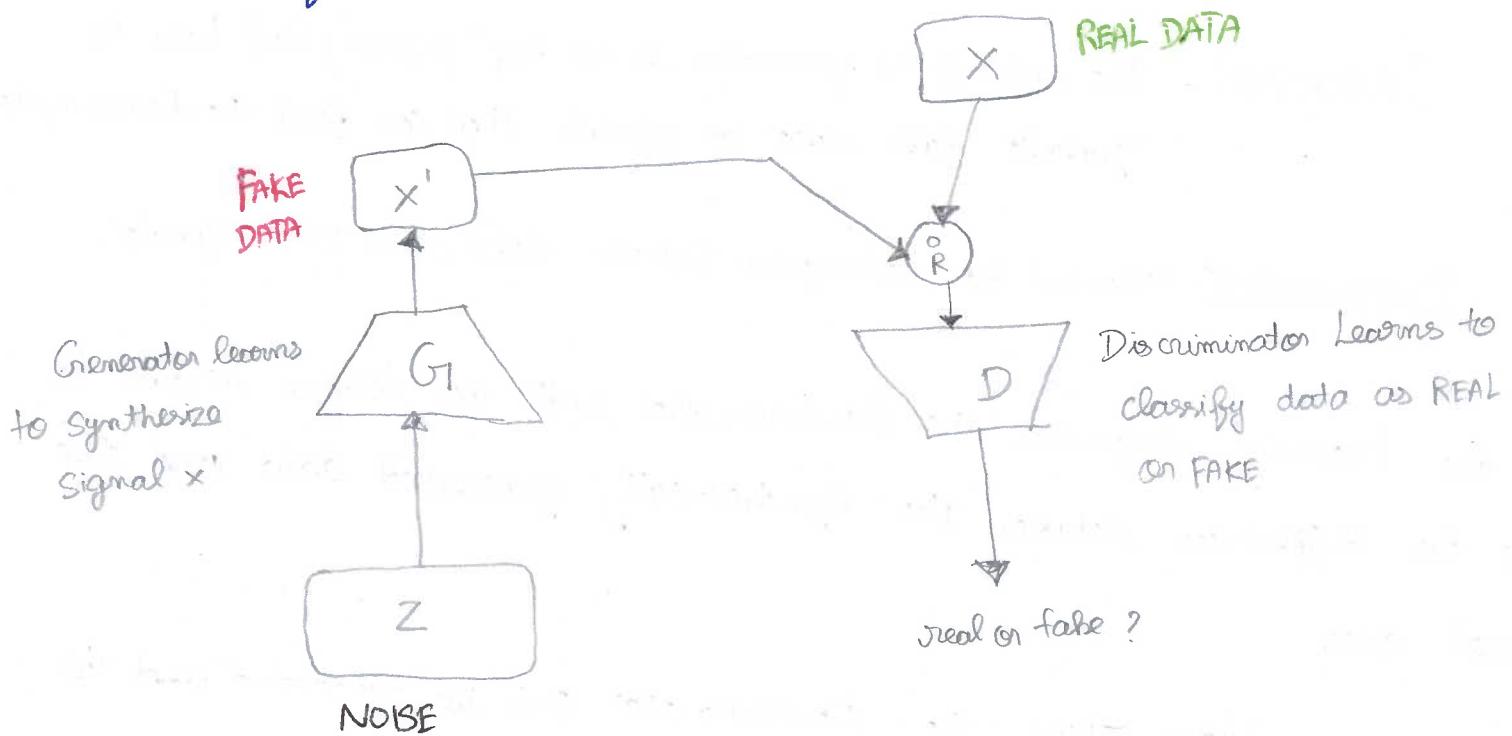
As the training progresses, the discriminator will no longer be able to see the differences between the synthetically generated data and the real ones.

From there, the discriminator can be **DISCARDED** and the generator can now be used to create new realistic signals that have never been observed before.

### CHALLENGE IN GAN:

- 1) There must be a **HEALTHY** competition between the **generator** and the **discriminator** in order for both networks to be able to learn simultaneously
- 2) **Loss Function** is computed from the output of the discriminator its parameter update is fast

- D) When the **discriminator** converges faster, the generator no longer receives sufficient GRADIENT UPDATES for its parameters and fails to converge.



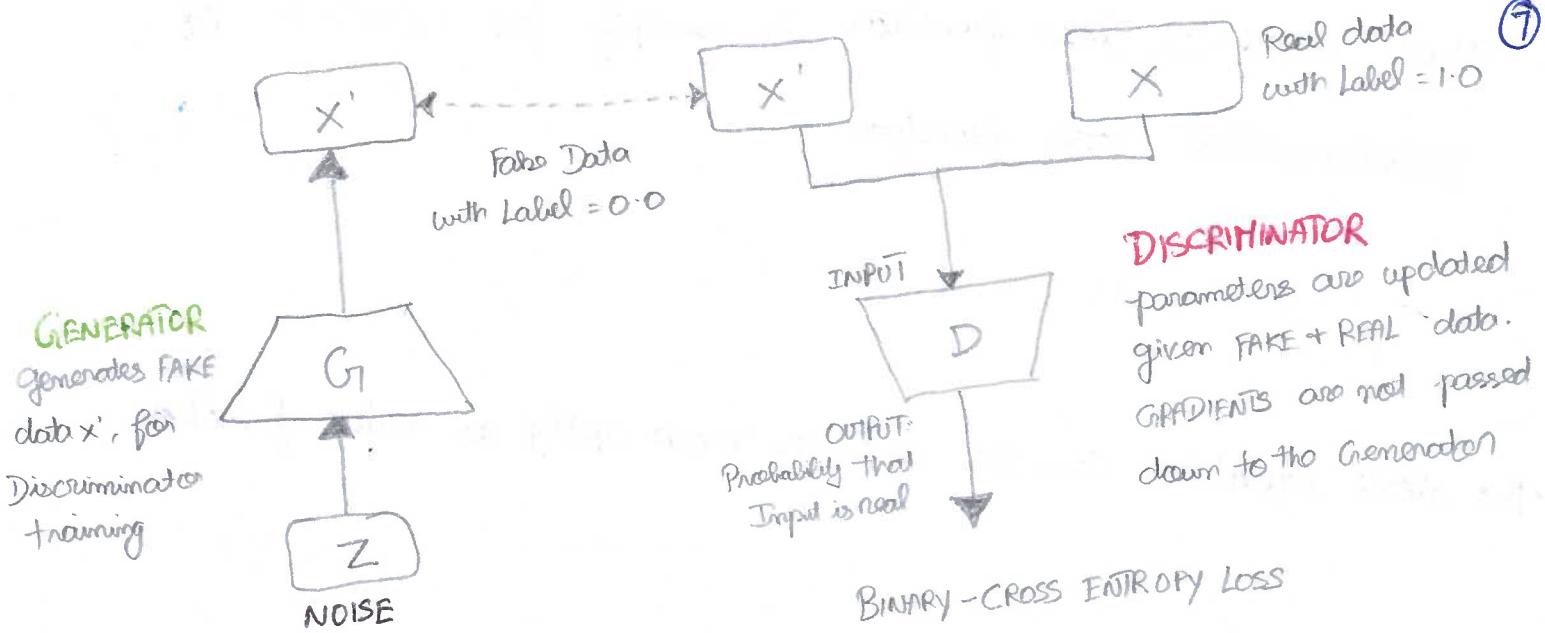
### PRINCIPLES OF GANs

- 1) **DISCRIMINATORS** are taught how to determine if a data is **REAL** or **FAKE**
- 2) **SAMPLES** of Real and fake data are used to train the **DISCRIMINATOR**
- 3) From time to time the **GENERATOR** will attempt to pretend that it generates **REAL** data
- 4) Initially, the **DISCRIMINATOR** will not be fooled and will tell the **GENERATOR** why the data is **FAKE**.

- 5) Taking into consideration this feedback, the GENERATOR hones his skills again and attempt to produce new fake data.
- 6) The DISCRIMINATOR again will be able to both spot the data as fake and justify why the data are fake.
- 7) As training progresses the GENERATOR masters his skills in making fake datas that are indistinguishable from real ones.
- 8) the GENERATOR can then infinitely print fake datas without getting caught by the DISCRIMINATOR as they are no longer identifiable as FAKE.

- \* Since the labeling process is automated, GANs are still considered part of the UNSUPERVISED LEARNING approach.
- \* The OBJECTIVE of the DISCRIMINATOR is to learn from this supplied dataset on how to distinguish REAL and FAKE data. During this part of GAN training, only the DISCRIMINATOR parameters will be updated.
- \* Like a typical binary classifier, the DISCRIMINATOR is trained to predict on a range of 0.0 to 1.0 in confidence values on how close a given input signal is to the True one.

- \* At regular intervals, the GENERATOR will pretend that its output is a genuine data and will ask the GAN to label it as 1.0
- \* The GENERATOR will use the gradients to update its parameters and improve its ability to synthesize fake signals.
- \* OVERALL, the whole process is akin to two networks COMPETING with one another while still CO-OPERATING at the same time.
- \* When the GAN training CONVERGES, the end result is a GENERATOR that can synthesize signals.
- \* the DISCRIMINATOR thinks these synthesized signal / data are real or with a label near 1.0, which means the DISCRIMINATOR can then be discarded.
- \* The GENERATOR part will be useful in producing meaningful outputs from arbitrary noise inputs.



As shown in the preceding figure, the **DISCRIMINATOR** can be trained by **MINIMIZING** the **LOSS FUNCTION**

$$L^{(D)}(\theta^{(G)}, \theta^{(D)}) = -E_{x \sim P_{\text{data}}} [\underbrace{\log D(x)}_{\text{DISCRIMINATOR OUTPUT FOR REAL DATA } x} - E_z \log (1 - D(G(z)))]$$

STANDARD BINARY CROSS-ENTROPY COST FUNCTION

- \* In order to minimize the loss function, the **DISCRIMINATOR** parameters  $\theta^{(D)}$  will be updated through **BACK-PROPAGATION** by correctly identifying the genuine data  $D(x)$  and synthetic data  $1 - D(G(z))$ .
- \* Correctly identifying real data is equivalent to  $D(x) \rightarrow 1.0$
- \* Correctly identifying fake data is equivalent to  $D(G(z)) \rightarrow 0.0$
- \*  $Z$  is the arbitrary encoding or noise vector that is used by **GENERATOR** to Synthesize new signal

- \* The GENERATOR loss function is simply the NEGATIVE OF DISCRIMINATOR loss function

$$L^{(G)}(\theta^{(G)}, \theta^{(D)}) = -L^{(D)}(\theta^{(G)}, \theta^{(D)})$$

The above equation can be re-written more aptly as value function

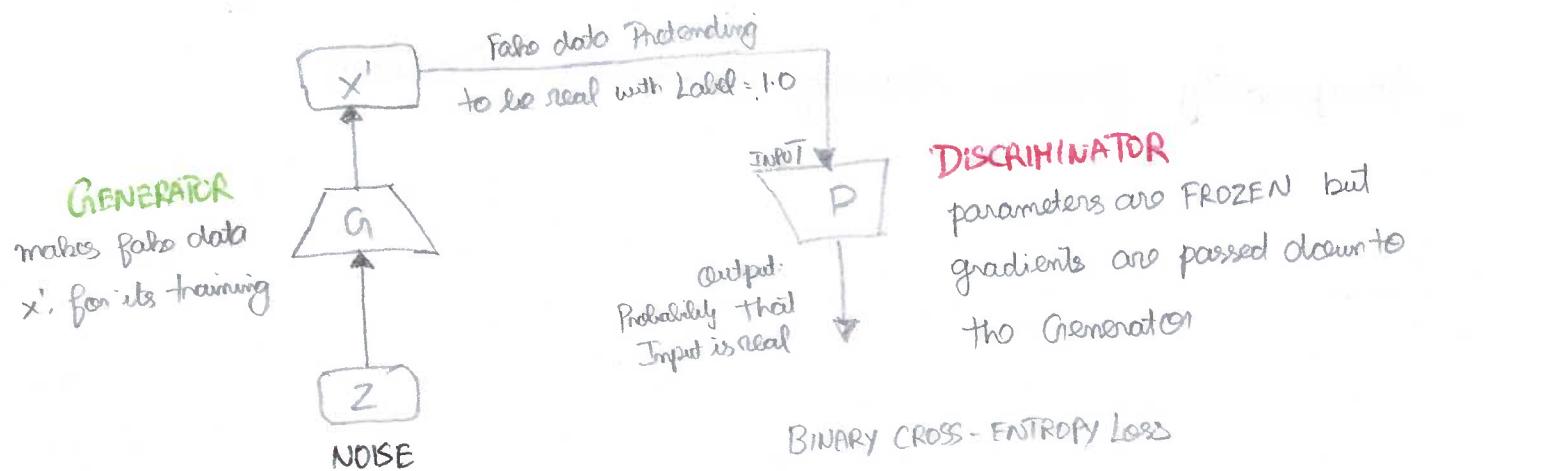
$$V^{(G)}(\theta^{(G)}, \theta^{(D)}) = -L^{(D)}(\theta^{(G)}, \theta^{(D)}) \quad \text{--- (1)}$$

- From the perspective of the GENERATOR, the above equation (1) should be MINIMIZED such that  $D(G(x))$  is close to 1 (Discriminator is fooled into thinking that  $G(x)$  is real)
- From the perspective of the DISCRIMINATOR, the above equation (1) should be MAXIMIZED such that  $D(x)$  is close to 1 and  $D(G(x))$  is close to 0
- Therefore, the generator training criterion can be written as a minimize problem

$$\theta^{(G)*} = \arg \min_{\theta^{(G)}} \min_{\theta^{(D)}} V^{(D)}(\theta^{(G)}, \theta^{(D)})$$

- occasionally, we will try to "Fool" the discriminator by pretending that the synthetic data is real with label 1.0

- 5) By MAXIMIZING with respect to  $\Theta^{(D)}$ , the optimizer  
Sends gradient updates to the DISCRIMINATOR parameters  
to consider this synthetic data as REAL.
- 6) At the same time by MINIMIZING with respect to  $\Theta^{(G)}$ , the  
optimizer will train the GENERATOR parameters on how  
to trick the DISCRIMINATOR.
- 7) However, in practice, the DISCRIMINATOR is confident in its  
prediction in classifying the Synthetic data as FAKE and  
will not update its parameters.
- 8) Furthermore, the gradient updates are small and have  
diminished significantly as they propagate to the GENERATOR  
layers
- 9) As a result, the GENERATOR fails to converge.



- 10) the solution is to reformulate the loss function of the GENERATOR

$$L^{(G)}(O^{(G)}, O^{(D)}) = -E_z \log D(G(z))$$

- 11) the above Loss Function Simply MAXIMIZES the chance of the DISCRIMINATOR into believing that the Synthetic data is real by training the GENERATOR.
- 12) the new formulation is no longer zero-SUM and is purely heuristics-driven.
- 13) In figure (Previous page), the GENERATOR parameters are only updated when the whole adversarial network is trained. this is because the gradients are passed down from the discriminator to the generator.
- 14) However, in practice the DISCRIMINATOR weights are only temporarily frozen during adversarial training.

# GAN TRAINING ALGORITHM :

for number of training iterations do

for K steps do

Label 0.0  $\leftarrow$  • Sample minibatch of "m" noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$

Label 1.0  $\leftarrow$  • Sample minibatch of "m" real examples  $\{x^{(1)}, \dots, x^{(m)}\}$

Discriminator ( $D_d$ ) { for a mini-batch, first the discriminator is trained on REAL and FAKE data }

• update the **DISCRIMINATOR** by ascending its stochastic gradient

DISCRIMINATOR NETWORK

$$\nabla_{D_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D_{d_i}(x^{(i)}) + \log (1 - D_{d_i}(G_{\theta_g}(z^{(i)}))) \right]$$

end for

{ new batch of FAKE Images }

Label 1.0  $\leftarrow$  • Sample minibatch of "m" noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$

• update the **GENERATOR** by ascending its stochastic gradient

ADVERSARIAL NETWORK

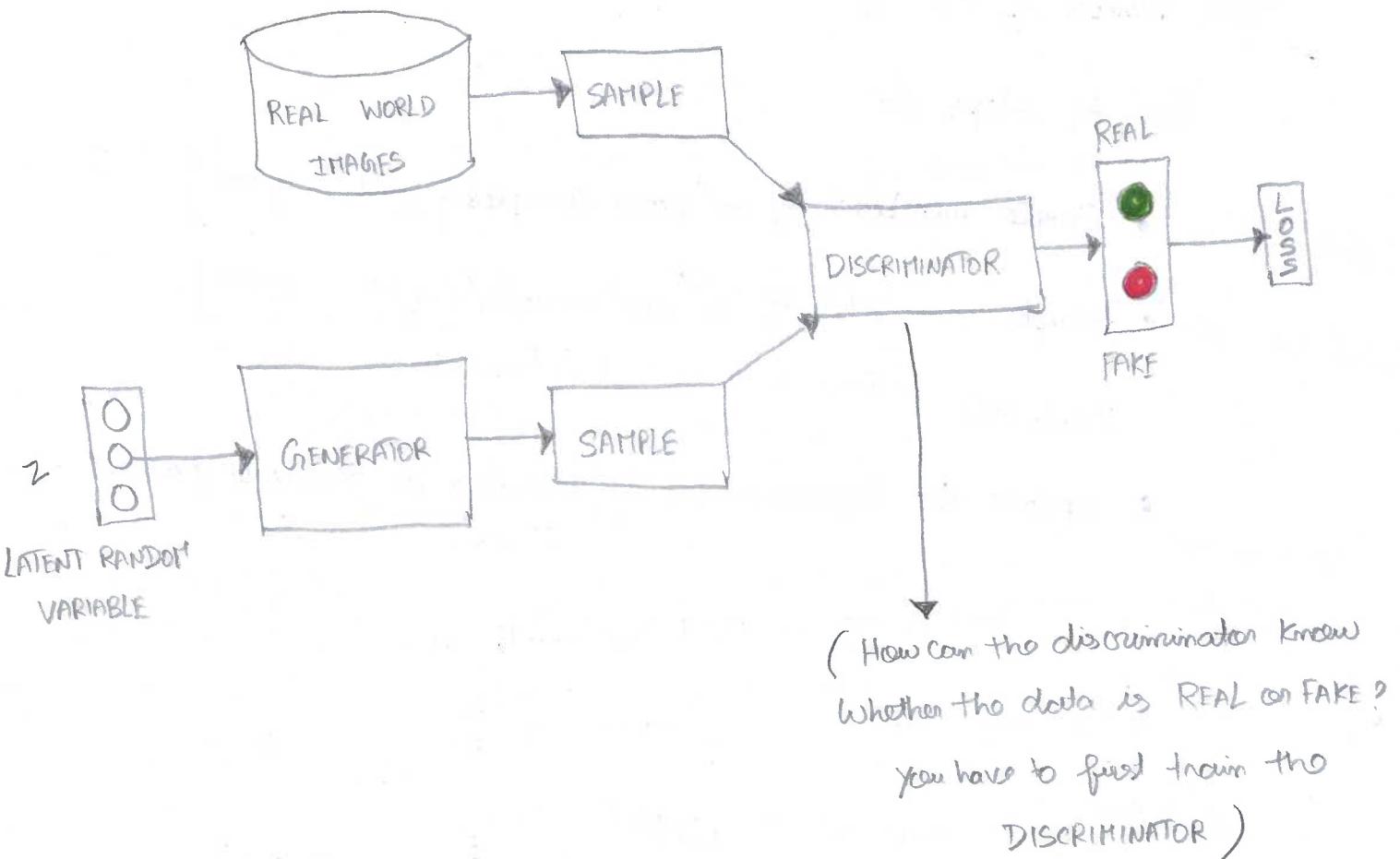
$$\nabla_{G_{\theta_g}} \frac{1}{m} \sum_{i=1}^m \log (D_{d_i}(G_{\theta_g}(z^{(i)})))$$

→ Here "DISCRIMINATOR" weights are frozen

end for

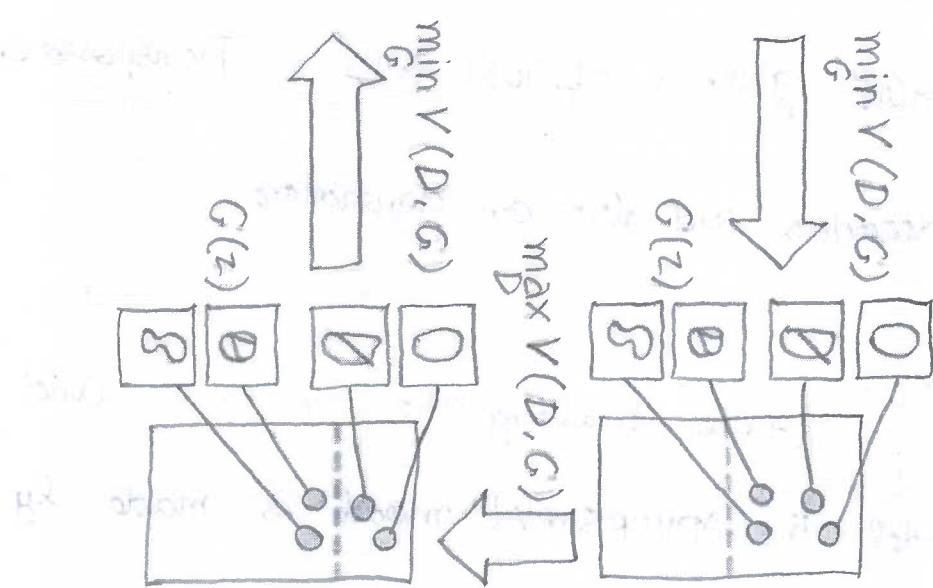
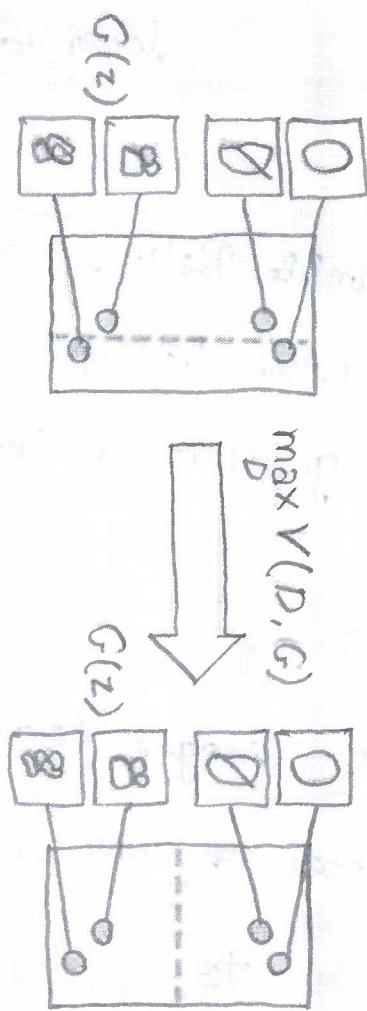
{ K → some people train the discriminator twice  
and  
train the generator once }

# ADVERSARIAL LEARNING



- 1) The Job of the "GENERATOR" is to make Images, Such that the "DISCRIMINATOR" Cannot recognize it as FAKE
- 2) "DISCRIMINATOR" train itself to be better at detecting FAKE Images
- 3) So "GENERATOR" and "DISCRIMINATOR" compete with each other
- 4) you train both of them at the Same time

$$\min_{G} \max_D V(D, G)$$



## BATCH NORMALIZATION:

Batch Normalization is used to stabilize learning by normalizing the input to each layer to have zero mean and unit variance.

Leaky ReLU: Unlike ReLU instead of zeroing out all outputs when the input is less than zero, Leaky ReLU generates a small gradient equal to  $\alpha \times \text{Input}$ .

## What is Conv2DTranspose?

(we can imagine transposed CNN (Conv2DTranspose) as the reversed process of CNN. In a simple example, if a CNN converts an IMAGE to FEATURE MAPS, a transposed CNN will produce an IMAGE given a FEATURE MAPS. Transposed CNN are used in decoders and also in generators.)

## ADVERSARIAL MODEL:

After building the GENERATOR and DISCRIMINATOR model, the ADVERSARIAL model is made by concatenating the generator and discriminator networks.

DCGAN (Deep convolutional Generative Adversarial Networks) :- (ii)

# build GENERATOR model

# build DISCRIMINATOR model

# build ADVERSARIAL model

# ADVERSARIAL = GENERATOR + DISCRIMINATOR

adversarial = Model(inputs, discriminator(generation(inputs)),  
name = model\_name)

Adversarial network training:-

Adversarial training is nothing but GENERATOR  
and DISCRIMINATOR training , with discriminator weights FROZEN

- \* The generator is trained via an adversarial network
- \* Alternatively train DISCRIMINATOR and ADVERSARIAL networks  
by batch. Discriminator is trained first with properly  
real and fake Images . Adversarial is trained next with  
fake Images pretending to be real.

**Min Max**  $V(D, G)$

$G \quad D$

$$V(D, G) = E_{x \sim P_{\text{data}}(x)} [\log D(x)] + E_{z \sim P_z(z)} [\log (1 - D(G(z)))]$$

$G$  = generator

$D$  = Discriminator

$P_{\text{data}}(x)$  = distribution of real data

$P(z)$  = distribution of generation

$x$  = Sample from  $P_{\text{data}}(x)$

$z$  = Sample from  $P(z)$

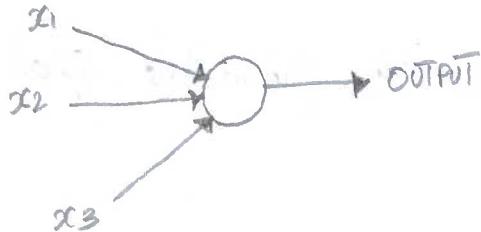
$D(x)$  = Discriminator network

$G(z)$  = Generation network

## How Perceptrons Work?

(12)

A perceptron takes several binary inputs  $x_1, x_2, \dots$  and produces a single binary output.

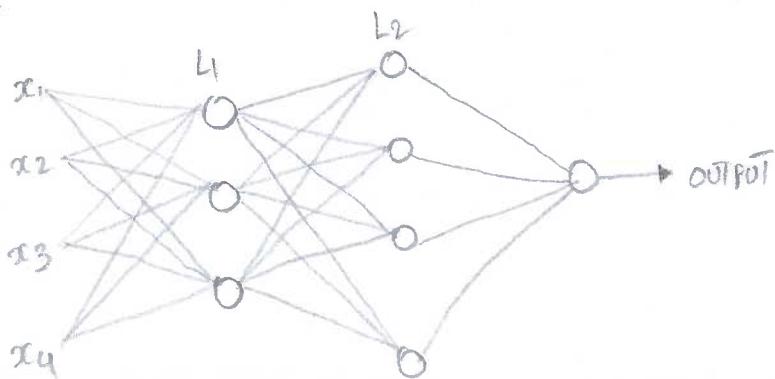


## What are weights?

Weights  $w_1, w_2, \dots, w_n$  are real numbers expressing the ~~respective~~  
**Importance** of the respective inputs to the output.

The neuron's output 0 or 1 is determined by whether the weighted sum  $\sum_j w_j x_j$  is less than or greater than some **threshold value**.

$$\text{output} = \begin{cases} 0 & \text{if } \sum_j w_j x_j \leq \text{threshold} \\ 1 & \text{if } \sum_j w_j x_j > \text{threshold} \end{cases}$$



In this network, the first column of perceptrons - first layer of perceptrons is making 3 very simple decision by weighing the Input evidence.

The Perceptrons in the 2<sup>nd</sup> Layer is making a decision by weighing up the results from the first layer of decision-making.

In this way, the perceptions in the 2<sup>nd</sup> Layer can make a decision at a more complex and abstract level than perceptions in the first layer.

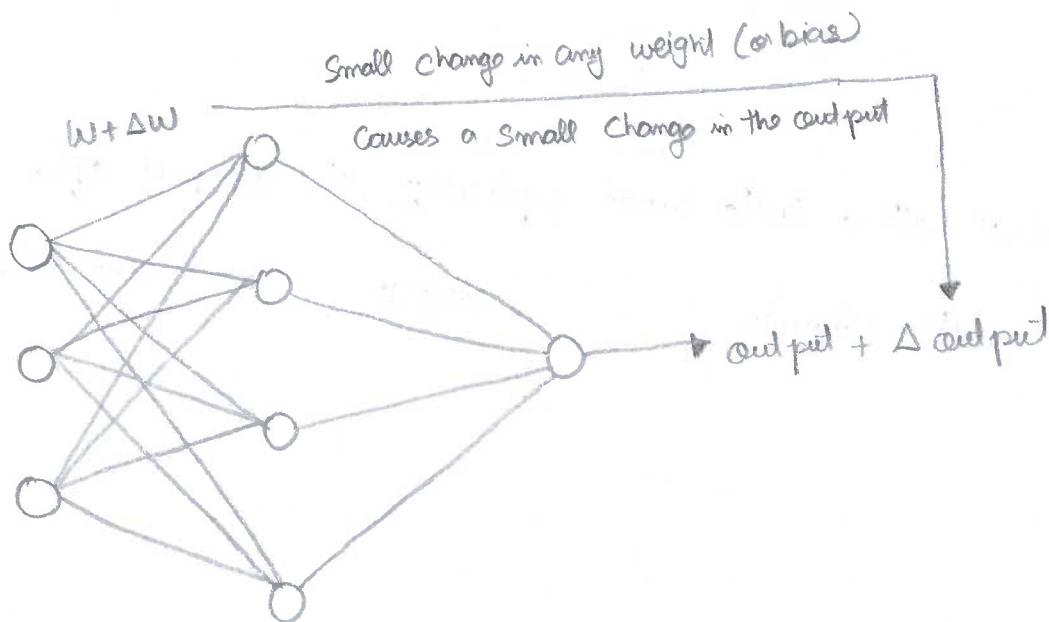
In this way, a many-layer network of perceptions can engage in sophisticated decision making

What are bias?

- i) you can think of bias as a measure of how easy it is to get the perceptron to output a 1 or, the bias is a measure of how easy it is to get the perceptron to fire

$$\text{Output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

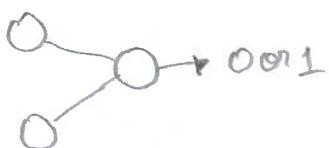
- \* It turns out that we can devise LEARNING ALGORITHM which can automatically tune the weights and biases of a network of artificial neurons (13)



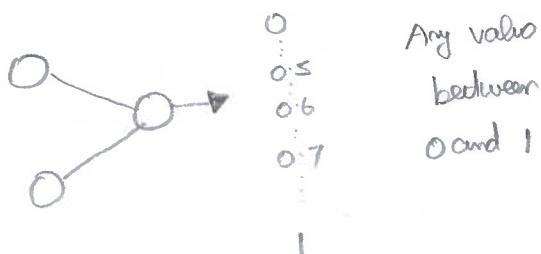
- \* If it is true that a small change in weight or bias causes only a small change in the output, then we could use this fact to modify the weights and biases to get our network to behave more in the manner we want.

What is the difference between Perceptron & Sigmoid neuron?

Perceptron output



Sigmoid neuron output



$\sigma(w \cdot x + b)$  where  $\sigma$  is called SIGMOID FUNCTION

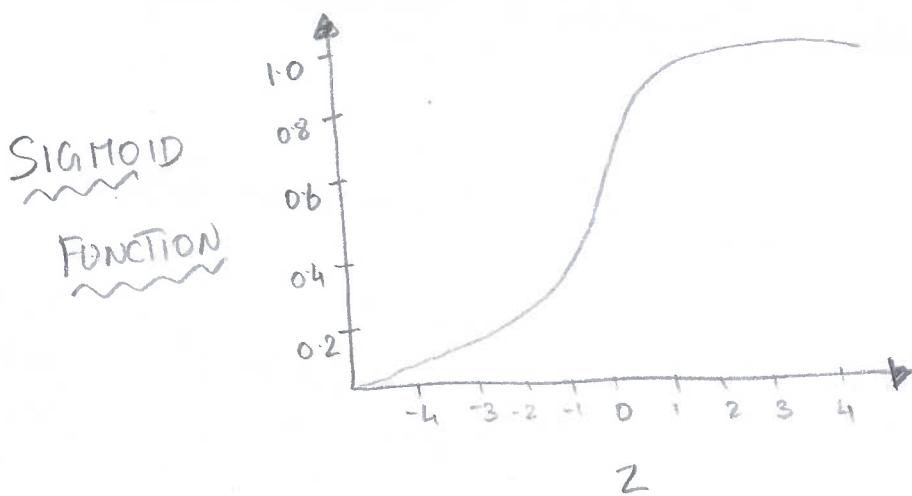
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

To put it all a little more explicitly, the output of a Sigmoid neuron with inputs  $x_1, x_2 \dots$ , weights  $w_1, w_2 \dots$  and bias  $b$  is

$$\frac{1}{1 + \exp(-\sum w_j x_j - b)}$$

$g_f$   
 $Z = w \cdot x + b$  is a Large Positive number  $\approx \sigma(z) \approx 1$

$g_b$   
 $z = w \cdot x + b$  is a very negative number  $\approx -\sigma(z) \approx 0$



$$\Delta \text{output} \approx \sum_j \frac{\partial \text{output}}{\partial w_j} \Delta w_j + \frac{\partial \text{output}}{\partial b} \Delta b$$

Change in weight                      Change in bias

Change in output is approximately equal to

What it derivative tell us is, weight ah change Parma, output also change aaguthu  
bias ah change Parma output also change aaguthu

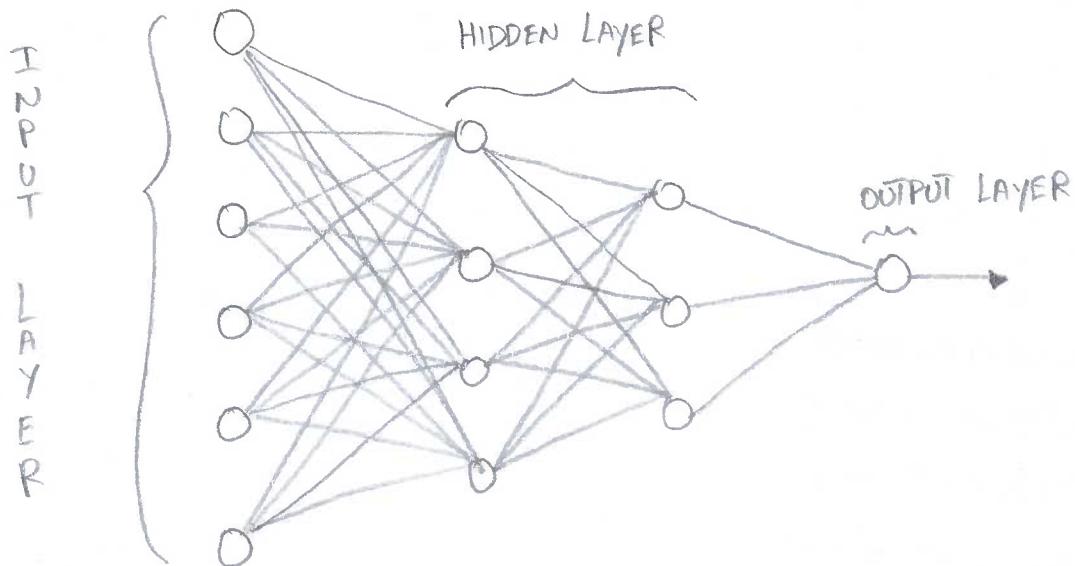
How should we interpret the output from a SIGMOID NEURON?

One big difference between the Perception and Sigmoid

neurons is the Sigmoid neurons don't just output 0 or 1. They can have as output any real number between 0 and 1, so values such as 0.173, 0.689, 0.04, 0.16... are legitimate outputs. This

can be useful, if we want to use the output value to represent the Intensity of the pixel in the Image to a neural network.

## MULTI-LAYER PERCEPTRONS



D) Multiple Layer networks are called MULTI LAYER PERCEPTRONS  
It's actually SIGMOID NEURONS

What is Feed Forward network?

The output from one layer is used as input to the next layer. Such networks are called "FEED FORWARD NETWORKS".

- \* this means there are **no loops** in the network
  - \* Information is always **Feed forward**

# NEURAL NETWORK CODE

(15)

The centerpiece is a NETWORK CLASS which we use to represent a neural network

Class Network (Object):

def \_\_init\_\_(self, sizes):

self.num\_layers = len(sizes)

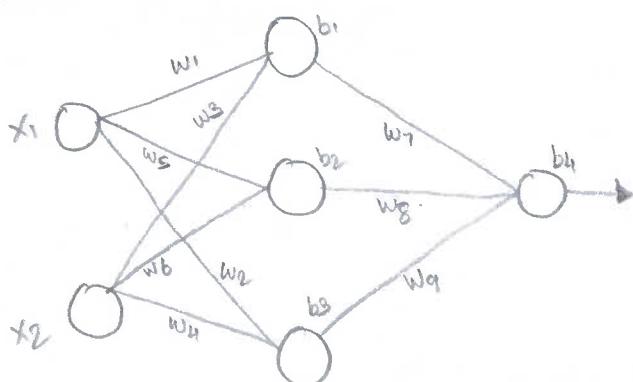
self.sizes = sizes

self.biases = [np.random.rand(y, 1) for y in sizes[1:]]

self.weights = [np.random.rand(y, x)

for x, y in zip(sizes[:-1], sizes[1:])] ]

- \* In Python3 Zip wraps two or more iterators with a Lazy generator.
- \* the zip generator yields tuples containing the next value from each iterator



net = Network ([2, 3, 1])

Input      no of neurons  
in 1<sup>st</sup> Layer

output

Sizes = [2, 3, 1]

GENERATE GAUSSIAN DISTRIBUTION  
WITH MEAN "0" AND SD "1"

Self-bias = [np.random.randint(y, 1) for y in sizes[1:]]

[3, 1]

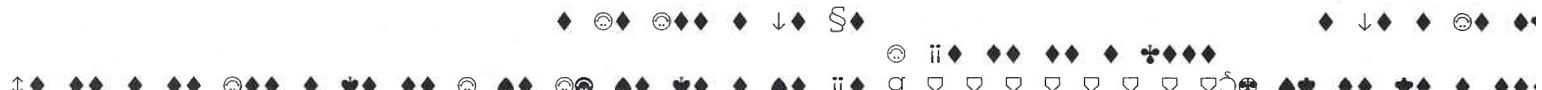
Because we don't need bias for the Input Layer

Returns List of numpy array

$$\text{Self-bias} = \left[ \begin{array}{c} \text{np.array}\left(\begin{bmatrix} [-0.8167], [-0.91592], [0.21430] \end{bmatrix}\right), \\ (3,1) \end{array} \right]$$

Self.weights =  $\left[ \text{np.random.randint}(y, x) \text{ for } x, y \text{ in} \right.$   
 $\left. \text{zip}(\text{size[:-1]}, \text{size[1:]}) \right]$

Returns List of numpy array



Self.weights =  $\left[ \text{np.random.randn}(y, x) \text{ for } x, y \text{ in } \text{Zip}(\text{sizes}[:i], \text{sizes}[i:]) \right]$

$\begin{matrix} [2, 3] & [3, 1] \\ x_1 \quad x_2 & y_1 \quad y_2 \end{matrix}$

Returns List of Numpy array

$\rightarrow$  weights connecting Input layer & the I<sup>st</sup> layer

Self.weights =  $\left[ \text{np.array} \left( \left[ \begin{matrix} [-0.90317, -0.79685], \\ [-1.0399, -1.13921], \\ [0.7169, 0.3896] \end{matrix} \right] \right), \right.$

$\rightarrow$  weights connecting I<sup>st</sup> layer & output layer

$\left. \text{np.array} \left( \left[ \begin{matrix} [-0.54459, -1.49075, -1.05287] \end{matrix} \right] \right) \right]$

### Sigmoid Function Code

```
def Sigmoid(z):
    return 1.0 / (1.0 + np.exp(-z))
```

\* we then add a Feed Forward method to the Network class, which given an input "a" for the network, returns the corresponding output

def feedforward(self, a):

for b, w in zip(self.biases, self.weights):

a = sigmoid(np.dot(w, a) + b)

return a

vector of activations of a layer

a is a  $(n, 1)$  numpy ndarray

where "n" is the number of

Input to the layer

\* of course, the main thing we want our Network object to do is

to learn. we will Implement SGD

def SGD(self, training\_data, epochs, mini\_batch\_size, eta, test\_data=None):

if test\_data: n-test = len(test\_data)

n = len(training\_data)

for j in range(0, epochs):

random.shuffle(training\_data)

`mini_batches = [ training_data[k:k + mini_batch_size]`

`for k in range(0, n, mini_batch_size)`

Returns List of List

`[ [ (x, y), (x, y), (x, y)],` → LIST OF TUPLES (Training data)

`[ (x, y), (x, y), (x, y)],`

`[ (x, y), (x, y), (x, y)],`

`[ (x, y), (x, y), (x, y)] ]`

`for mini_batch in mini_batches:`

`self.update_mini_batch(mini_batch, eta) →`

UPDATES NETWORK  
WEIGHTS AND  
BIASES

`If test_data:`

`Print "Epoch {0}: {1} | {2}" . format(j, self.evaluate(test_data), n-test)`

`else:`

`Print "Epoch {0} complete". format(j)`

def update\_mini\_batch (self, mini-batch, eta):

nabla\_b = [np.zeros(b.shape) for b in self.biases]

nabla\_w = [np.zeros(w.shape) for w in self.weights]

for x, y in mini-batch:

delta\_nabla\_b, delta\_nabla\_w = self.backprop(x, y)

nabla\_b = [nb + dnb for nb, dnb in zip(nabla\_b, delta\_nabla\_b)]

nabla\_w = [nw + dnw for nw, dnw in zip(nabla\_w, delta\_nabla\_w)]

Self.weights = [w - (eta / len(mini-batch)) \* nw

for w, nw in zip(self.weights, nabla\_w)]

Self.biases = [b - (eta / len(mini-batch)) \* nb

for b, nb in zip(self.biases, nabla\_b)]

\* update\_mini\_batch works simply by computing the GRADIENTS for every training example in the mini-batch and the updating

Self.weights

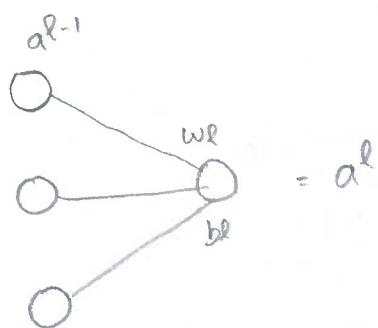
Self.biases

appropriately

## BACK PROPAGATION

$$a^l = \sigma(w^l a^{l-1} + b^l)$$

- \* this expression gives us a much more global way of thinking about how the activation in one layer relate to activations in the previous layer.



- \* The **GOAL** of back propagation is to compute the partial derivatives of  $\frac{\partial C}{\partial w}$  and  $\frac{\partial C}{\partial b}$  of the Cost function  $C$  with respect to any weight "w" or bias "b" in the network.
- \* Back propagation is about understanding how changing the weights and bias in a network changes the Cost function
- \* Chain rule is used to compute the gradient

Why is learning slow?

Our neurons learns by changing the weight and the bias at a rate determined by the partial derivatives of the cost function  $\frac{\partial C}{\partial w}$  and  $\frac{\partial C}{\partial b}$ . So saying "Learning is slow" is really the same as saying that those partial derivatives are small.

### CROSS-ENTROPY COST FUNCTION:-

$$a = \sigma(z)$$

$$C = -\frac{1}{m} \sum_{i=1}^n \left[ y_i \ln a + (1-y_i) \ln (1-a) \right] \quad \text{--- (1)}$$

n = total number of items of Training data

### Properties

- \* The cost function is NON-NEGATIVE i.e ( $C \geq 0$ )
- \* Both Logarithm are of numbers in the range 0 to 1
- \* There is a minus sign out the front of the sum

- \* If the neuron's actual output is close to the desired output for all training inputs " $x$ ", then the cross-entropy will be close to "zero".

Example:

Suppose

$$y=0 \text{ and } a \approx 0 \text{ for some Input } "x"$$

→ This is the case when the neuron is doing a good job on that input.

→ First term in equation ① vanishes

How cross-entropy differs from quadratic cost?

Cross-entropy avoids the problem of

Learning Slowing Down

- \* The cross-entropy is **positive**, and tends toward zero as the neuron gets better at computing the desired output " $y$ " for all training inputs " $x$ ".

## Max-Pooling Intuition:

- \* we can think of max-pooling as a way for the network to ask whether a given feature is found anywhere in a region of the image.
  - It then throws away the exact positional information, L2 Pooling is a way of **CONDENSING** information from the **CONVOLUTIONAL** layer.

## How NETWORK PARAMETERS ARE UPDATED:

- \* The parameters of the network, filters in the convolutional layer, weight-matrices in the fully-connected layer and biases are **TRAINED** by **BACKPROPAGATING** the **DERIVATIVE** of the **Loss** with respect to the parameters throughout the network and updating the parameters via **STOCHASTIC GRADIENT DESCENT**.
- \* Understanding the operation of a CONVNET requires interpreting the feature activity in Intermediate Layers.
- \* The CONVNET uses **RELU** non-linearities, which rectify the feature maps, thus ensuring the feature maps are always **POSITIVE**.

## FEATURE VISUALIZATION:

- \* The visualizations reveal that a set of particular feature maps focuses on the "GRASS" in the background, not the foreground objects.

What causes high activations in a feature map?

- \* Layer 2 responds to corners and edges
- \* Layer 3 has more complex invariances like mesh patterns
- \* Layer 4 captures class-specific like dog's face, bird's leg
- \* Layer 5 shows entire object with significant pose variation

## ARCHITECTURE SELECTION:

- \* While visualization of a trained model give insight into its operation, it can also assist with selecting good architectures in the first place.
- \* The first and 2<sup>nd</sup> layer of Krizhevsky architecture, the 1<sup>st</sup> layer filters are a mix of extremely HIGH and LOW frequency information with the little coverage of MID frequencies

Does the Model learns the location of the object?

With "IMAGE CLASSIFICATION" approaches, a natural question is If the MODEL is truly identifying the LOCATION of the object in the IMAGE or Just using the Surrounding context.

This can be clarified by Systematically OCCLUDING different portions of the Input Image with a GREY SQUARE and monitoring the output of the CLASSIFIER.

The example clearly shows the MODEL is LOCALIZING the objects with the SCENE, as the probability of the correct class drops significantly when the object is OCCLUDED.

## WHY CONVOLUTION ?

The first thing to have in mind is that an Input Volume is normally Convolved with several different KERNELS

Why?

Because, we expect each different Kernel to extract different features from the Input Image.

CONVNETS are Inherently translation Invariant.

## QuickNAT

### ABSTRACT

- \* Pre-train on auxiliary labels created from existing Segmentation Software
- \* Pre-trained model is fine TUNED on MANUAL LABELS to rectify errors in auxiliary labels.

### INTRODUCTION

- \* In, order to access measurements like VOLUME, THICKNESS OR SHAPE OF A STRUCTURE, the neuroanatomy needs to be SEGMENTED
- \* Existing ATLAS-BASED methods require hours of processing time for each Scan and may result in SUB-OPTIMAL Solutions.
- \* the Limited availability of training data with Manual annotations present the main challenge in extending F-CNN models to brain Segmentation

D) AUXILIARY LABELS - Freesurfer generated Segmentation

2) MANUAL LABELS - Expert clinician Segmentation

- \* Pre-training provides a good prior initialization of the network

## METHODS

- D Given an Input MRI brain Scan  $I$ , we want to Infer its segmentation map  $S$ , which indicates **27 CORTICAL** and **SUB-CORTICAL** structures.

## ARCHITECTURE:

- D QuickNAT has an **ENCODER / DECODER** like 2D F-CNN architecture with 4 encoders and 4 decoders separated by a **BOTTLE NECK** layer.
- 2) Final Layer is a classifier block with "SOFTMAX"
- 3) The architecture includes **SKIP CONNECTIONS** between all encoder and decoder blocks, of the same spatial resolution
- 4) These skip connections not only provide **ENCODER FEATURE INFORMATION** to the **DECODER**, but also provide a path of gradient flow from the shallow layers to deep layers, improving training

## WHAT ARE SKIP-CONNECTIONS ?

(22)

Skip connections add encoder features to the decoders from aiding segmentation and thus allow gradients to flow from deeper to shallower regions of the network.

## CLASSIFIER BLOCK ?

The output feature map from the last decoder block is passed to the classifier block, which is basically a convolutional layer with  $1 \times 1$  kernel size that maps the input to a  $N$  channel feature map, where " $N$ " is the number of classes (28 in our case). This is followed by a SOFTMAX LAYER to map the activations to PROBABILITIES

- \* In Probability theory, the output of the SOFTMAX FUNCTION can be used to represent a categorical distribution - that is, a probability distribution over " $K$ " different possible outcomes.

## NON-LINEAR FUNCTION

- \* The main terminologies needed to understand for non-linear functions are

Derivative } Change in y-axis w.r.t Change in x-axis. It is  
Differential } also known as slope

Monotonic } A function which is either non-increasing  
function } or non-decreasing

## SIGMOID OR LOGISTIC ACTIVATION FUNCTION

- 1) The Sigmoid function curve looks like an S-shape
- 2) The main reason why we use sigmoid function is because it exists between (0 to 1). Therefore, it is especially used for models where we have to "predict the probability" as an output.
- 3) The function is differentiable. That means, we can find the slope of the sigmoid curve at any two points
- 4) The function is monotonic but function's derivative is not.

## BOTTLE NECK LAYER:

(23)

- 1) This layer forces the ENCODER to Compress the Input Image and create some more high level representation
- 2) The DECODER will also have to learn how to decompress the high level representation into something that looks like an Input Image.
- 3) we have to just switch what "output target is" from trying to recreate the Input Image into trying to create the **SEGMENTATION** maps.

## FULLY CONVOLUTIONAL NETWORK - (U-Net)

- 1) Domain specific Image **FEATURES** are learnt using **SGD**
- 2) Learnt kernels are shared across all pixels
- 3) Image **CONVOLUTION** operation exploit the structural information in medical Images
- 4) Coarse feature-maps capture **CONTEXTUAL** Information and highlight the category and location of foreground objects

- 5) Feature-maps extracted at multiple scales are later MERGED through SKIP CONNECTIONS TO COMBINE Coarse and fine-level dense predictions.
- 6) Convolutional layer learns filters capturing local SPATIAL pattern along all the input channels and generates feature maps jointly encoding the SPATIAL and CHANNEL information.

SPATIAL → Means anything related to SPACE

CHANNEL → 3 for RGB Images

1 for GRAY Scale Images

SPATIAL INVARIANCE → Most of us can recognize specific objects under variety of conditions because we learn abstractions. These abstractions are thus invariant to

1) SIZE

2) CONTRAST

3) ROTATION

4) ORIENTATION

5) LOCATION

CONVNET

Successfully Captures  
the

SPATIAL  
and

TEMPORAL

Information

## SPATIAL

1) When you have a series of Images

When you are analyzing Image.

It includes Coordinates, Intensity, gradient, resolution, size, contrast, Rotation, orientation & Location

## TEMPORAL

1) When you have a series of Images taken at different time. Correlations between the Images are often used to monitor the dynamic change of the object

## What is SE all about?

Attention Gates Network (SE Block) explores the "SPATIAL" and "CHANNEL-WISE" patterns independently. Modelling the Interdependencies between the channels of feature maps.

2) Spatial dependency is a process in which the variables like gradient, Intensity, Size, Contrast, orientation + location may be dependent

