

16/08/2023

## Deep Learning

Advanced learning algorithms:

1). Neural Network 

- inference (prediction)
- training

2). Practical advice for building machine learning systems 

3). Decision Tree 

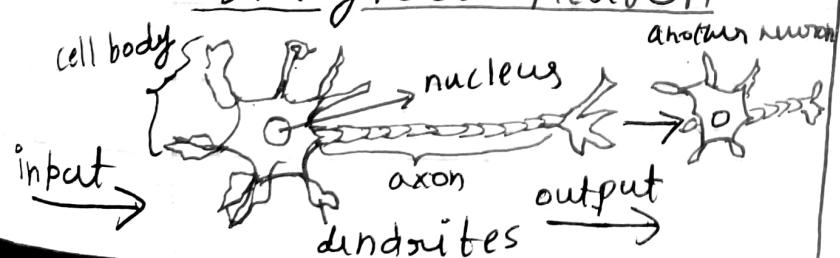
Neural networks intuition

(i) Neurons and the brain

### Neural Network

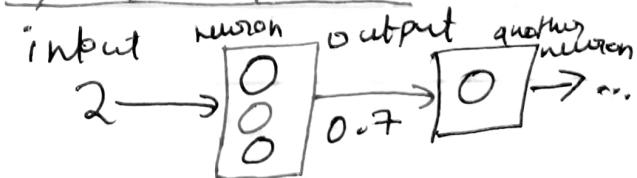
- origins: Algorithms that try to mimic the human brain.
- Used in the 1980's & early 1990's  
Fell out of favor in the late 1990's
- Resurgence from around 2005
- Speech → Images → text (NLP) → ...

### Biological neuron



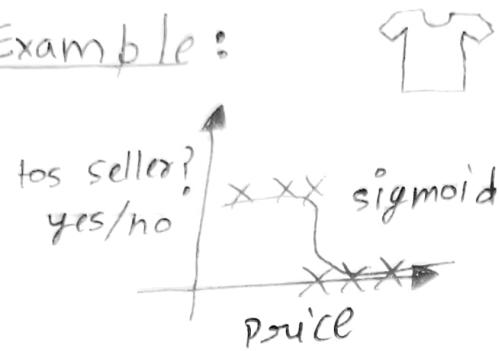
simplified mathematical model

of a neuron



## (ii) Demand Prediction

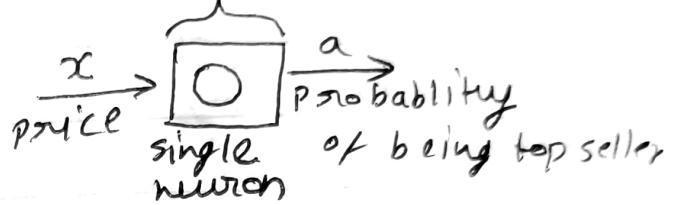
Example:



$x = \text{price}$  (Input)

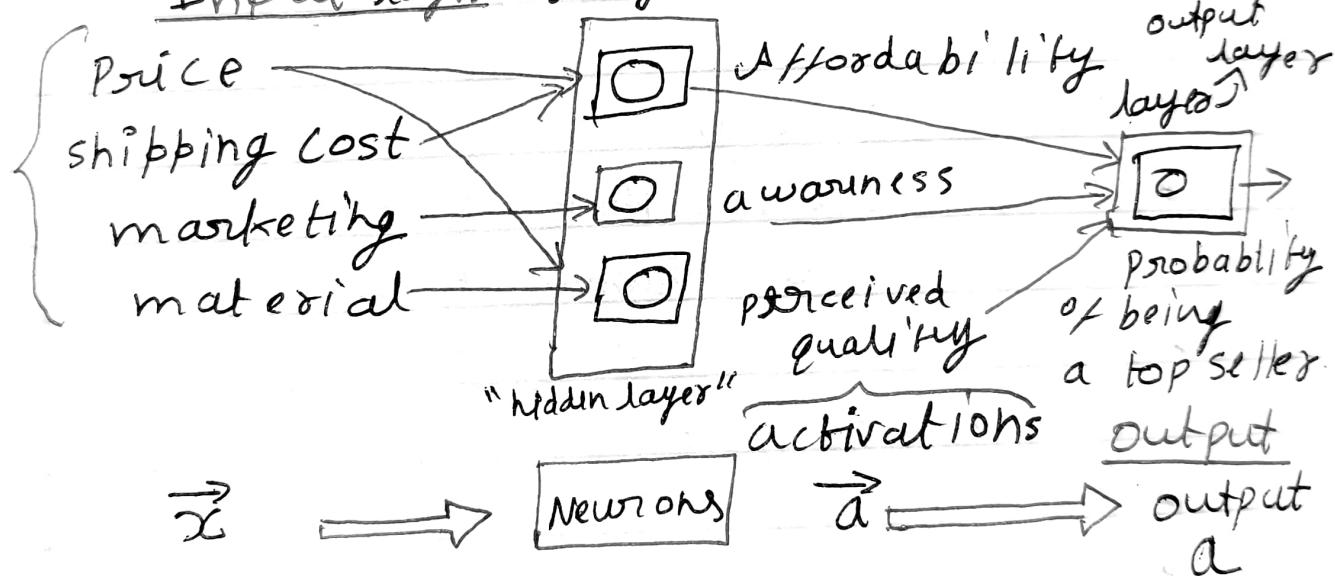
$$a = f(x) = \frac{1}{1 + e^{-(wx + b)}} \quad (\text{Output})$$

activation

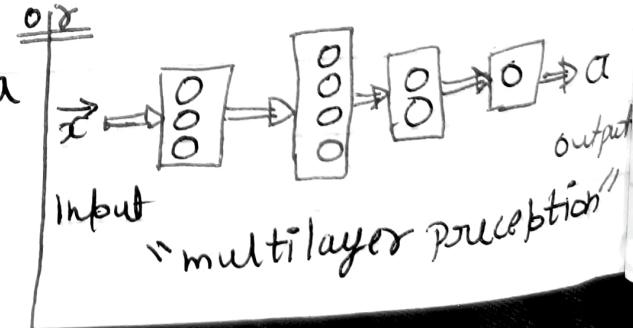
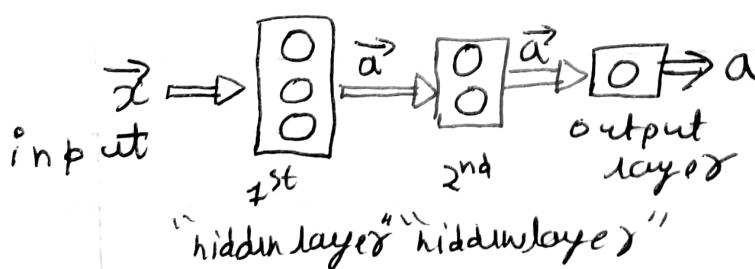


Example: if we have multiple features

Input layer



multiple Hidden layers



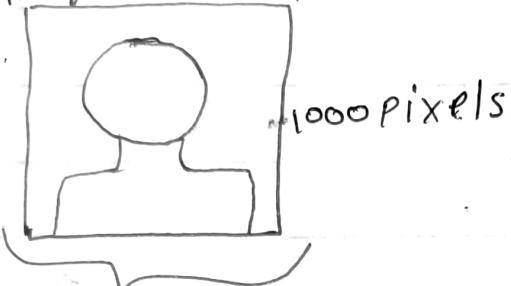
"multilayer perceptron"

(iii)

## Example: Recognizing Images

### Face recognition

Image 1000 pixels



197	185	203	000	000
000	57	64	92	000
0	8			
	000	187	214	

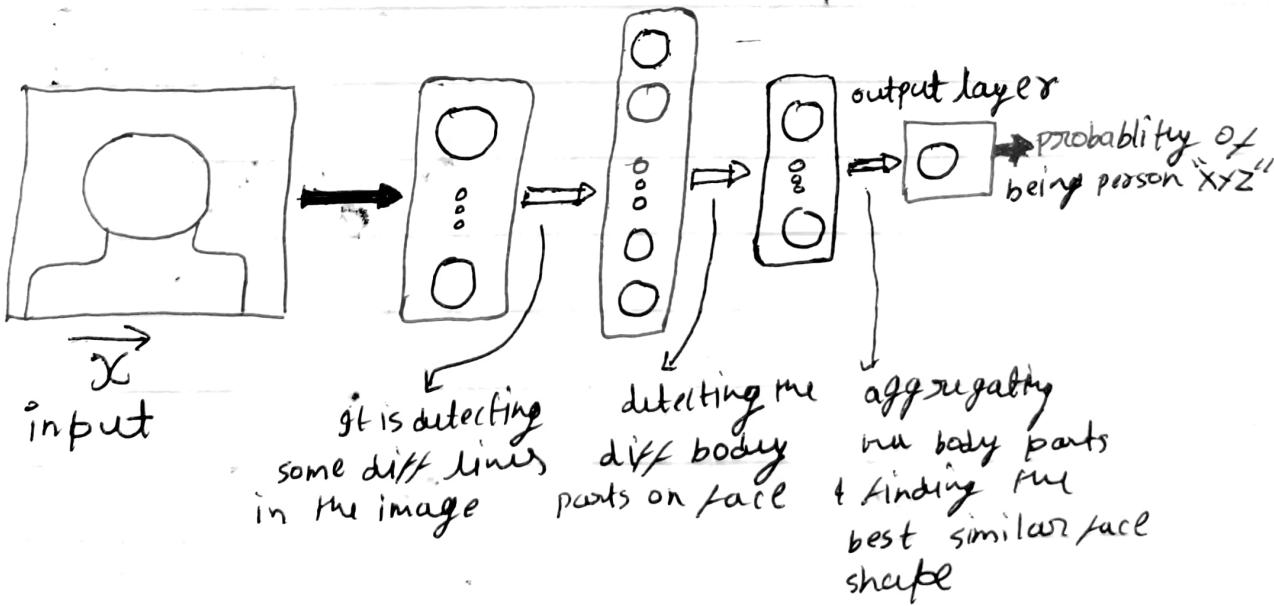
1000 Rows

1000 columns

$$\Rightarrow \vec{x} =$$

197
185
203
0
57
64
92
0
187
214

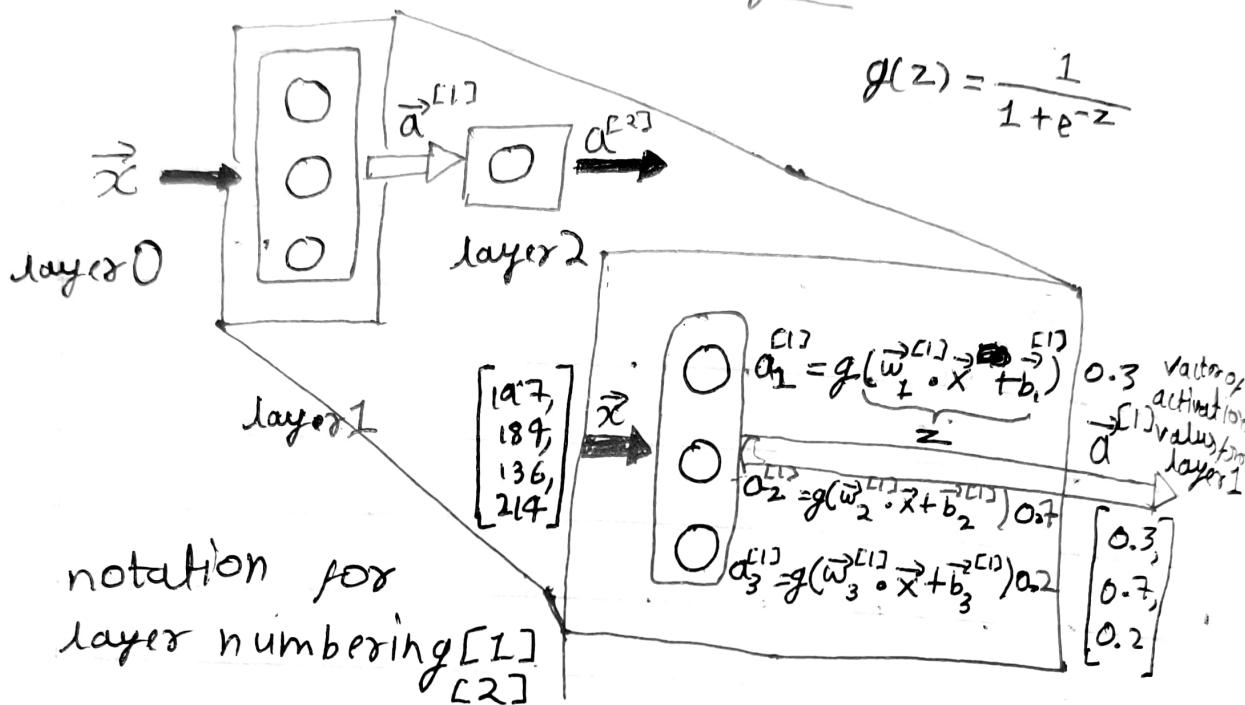
1 million



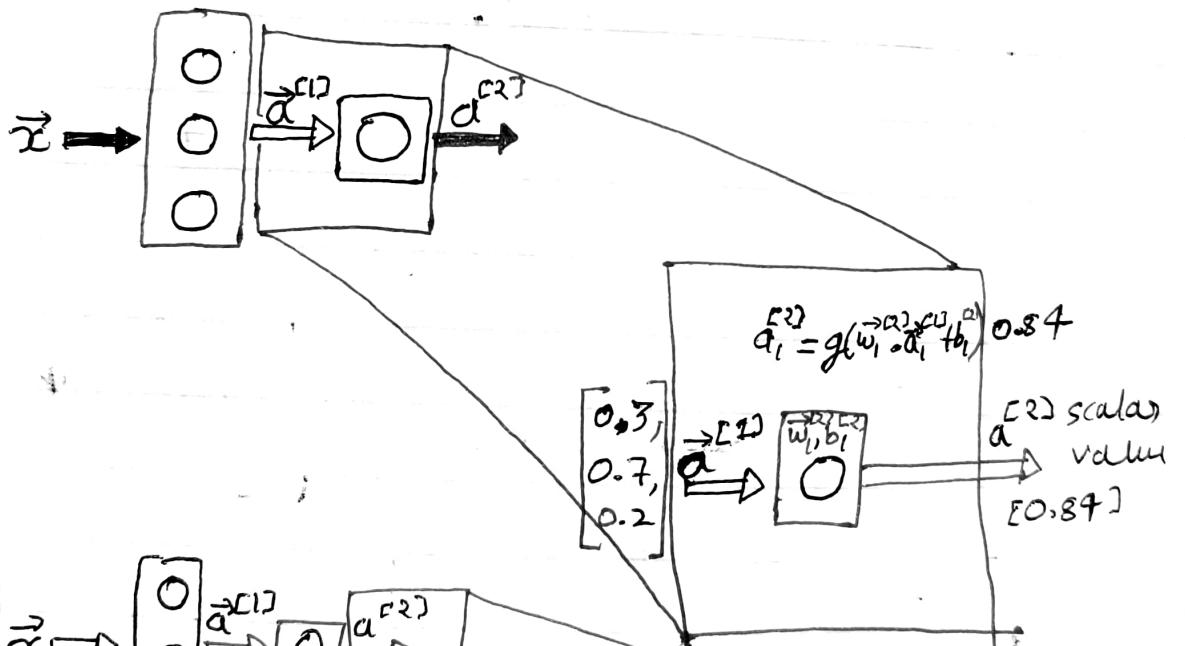
# Neural Network Model

## (i) Neural network layer

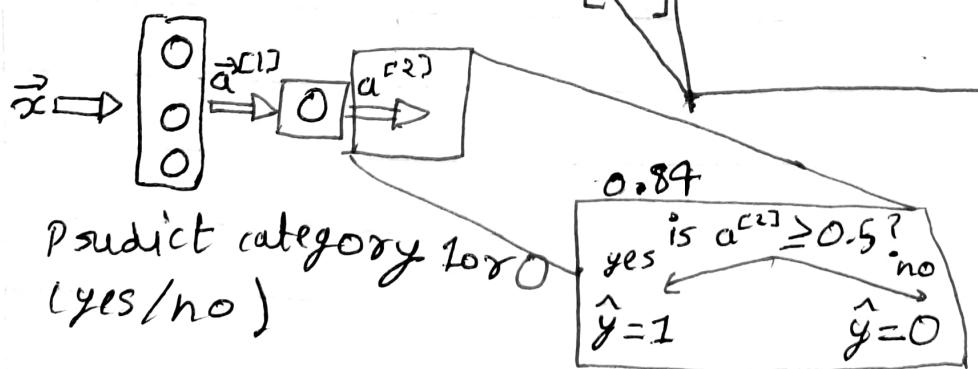
①



②



③



"activation function"

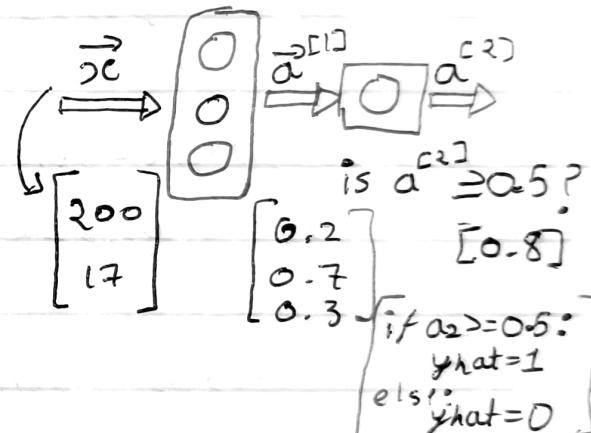
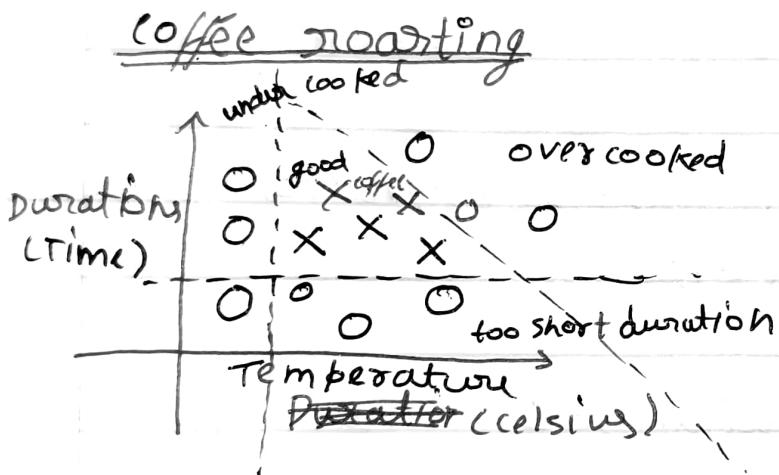
$$\Rightarrow a_j^{[l]} = g(w_j^{[l]} \cdot a^{[l-1]} + b_j^{[l]})$$

activation value of  
layer  $\ell$ , unit  $j$  (neuron)

Parameters  $w$  &  $b$   
of layer  $\ell$ , unit  $j$

## TensorFlow implementation

### (i) Inference in Code



```
⇒ x = np.array([200.0, 17.0])
```

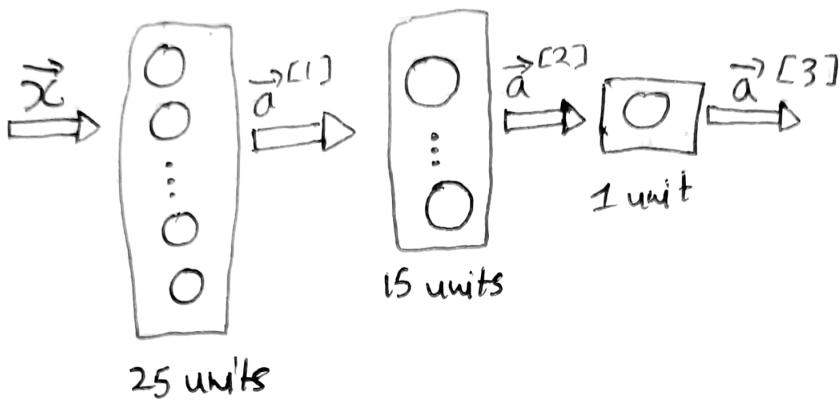
```
layer_1 = Dense(3, activation='sigmoid')
```

```
a_1 = layer_1(x)
```

```
⇒ layer_2 = Dense(1, activation='sigmoid')
```

```
a_2 = layer_2(a_1)
```

## model for digit classification



```
x = np.array([0.0, ... 245, ... 250, ... 0])
```

```
layer_1 = Dense(units=25, activation='sigmoid')
```

```
a1 = layer_1(x)
```

```
layer_2 = Dense(units=15, activation='sigmoid')
```

```
a2 = layer_2(a1)
```

```
layer_3 = Dense(units=1, activation='sigmoid')
```

```
a3 = layer_3(a2)
```

```
if a3 ≥ 0.5:
```

```
yhat = 1
```

```
else:
```

```
yhat = 0
```

inference code

## (ii) Data in TensorFlow

### Feature vectors

temperature (celsius)	duration (minutes)	Good coffee? (1/0)
200.0	17.0	1
425.0	18.5	0
...	...	...

$x = \text{np.array}([200.0, 17.0])$

$\boxed{[200.0, 17.0]}$

→ why double braces?

$\Rightarrow \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3}$  if we want to implement this in code then we have to use above format with double braces

$\Rightarrow x = \text{np.array}([ [1, 2, 3], [4, 5, 6] ])$

$\boxed{[ [1, 2, 3], [4, 5, 6] ]}$

\* Now when we used  $[[200.0, 17.0]]$  by double braces because we have to pass a 2-D array instead of 1-D array.

2-D vector {  
 $x = np.array([200.0, 17.0]) \rightarrow [200, 17]_{1 \times 2}$   
 $x = np.array([200.0], [17.0]) \rightarrow [200] \rightarrow [17]_{2 \times 1}$

1-D vectors { x = np.array ([200.0, 17.0]) ~~200.0, 17.0~~  
no rows & no cols

Going back to the code - ,  
`x = np.array([[200, 17]])`

```
x=np.array([[200,17]])
```

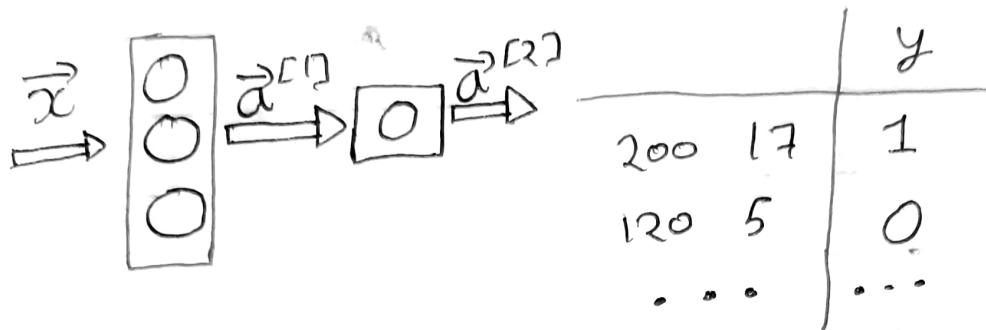
layer\_1=Dense(units=3, activation='sigmoid')

$$a_1^i = \text{layer}_1(x)$$

data in tensorflow {  
    } ⇒ print(a1)  
    [tf.Tensor([0.2, 0.7 0.3]), shape=(1,3), dtype=float]

data in Numpy  $\Rightarrow$  a.numpy()  
[array([0.2, 0.7, 0.3]), dtype= float32)]

### (iii) Building a neural Network



`layer-1 = Dense(units=3, activation='sigmoid')`  
`layer-2 = Dense(units=1, activation='sigmoid')`  
`model = Sequential([layer-1, layer-2])`

`x = np.array([ [200.0, 17.0],  
[120.0, 5.0],  
[425.0, 20.0],  
[212.0, 18.0] ])`  $4 \times 2$

`y = np.array([1, 0, 0, 1])`

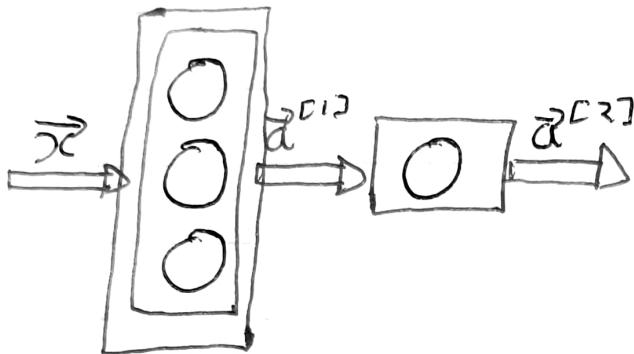
`model.compile(...)` ← more about this next week!  
`model.fit(x, y)`

`model.predict(x-new)`

we can do this also { `model = Sequential ([`  
`Dense(units=3, activation='sigmoid'),`  
`Dense(units=1, activation='sigmoid')]` }

# Neural network implementation in python

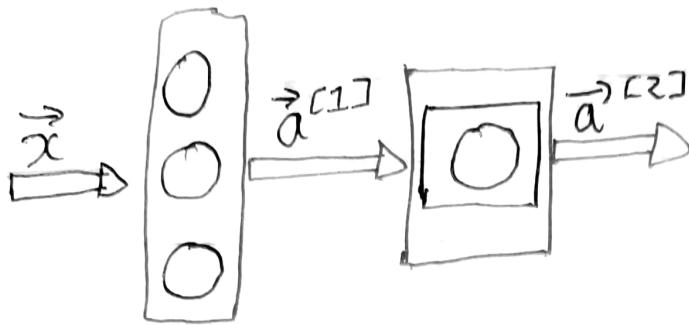
## (i) Forward prop in a single layer



```
x = np.array([200, 17])           # D array
#  $a_1^{[1]} = g(w_1^{[1]} \cdot \vec{x} + b_1^{[1]})$       #  $a_2^{[1]}$       #  $a_3^{[1]}$ 
```

$w_{1-1} = np.array([1, 2])$	$w_{1-2} = [-3, 4]$
$b_{1-1} = np.array([-1])$	$b_{1-2} = [1]$
$z_{1-1} = np.dot(w_{1-1}, x) + b_{1-1}$	$z_{1-2} = np.dot(...)$
$a_{1-1} = sigmoid(z_{1-1})$	$a_{1-2} = sigmoid(z_{1-2})$

$a_1 = np.array([a_{1-1}, a_{1-2}, a_{1-3}])$



$$a_1 = \text{np.array}([a_{1-1}, a_{1-2}, a_{1-3}])$$

$$\# a_2^{[2]} = g(\vec{w}_2^{[2]} \cdot \vec{a}^{[1]} + b_2^{[2]})$$

$$w_2-1 = \text{np.array}([-7, 8, 9])$$

$$b_2-1 = \text{np.array}[3]$$

$$z_2-1 = \text{np.dot}(w_2-1 \cdot a_1) + b_2-1$$

$$a_2-1 = \text{sigmoid}(z_2-1)$$

↓

$$a_2 = \text{np.array}[a_2-1]$$

(ii) General implementation of forward propagation

forward prop in Numpy

activation vector

```
def dense(a-in, W, b):  
    units = W.shape[1] # no. of columns  
    a-out = np.zeros(units) # [0,0,...,0]  
    for j in range(units):  
        w = W[:, j] # columns 1 by 1  
        z = np.dot(w, a-in) + b[j]  
        a-out[j] = g(z) # sigmoid fn  
    return a-out
```

~~def sequential(a, W, b, l):~~

~~for i in range(l):~~  
~~a = dense(x,~~

```
def sequential(X):  
    a1 = dense(X, w1, b1)  
    a2 = dense(a1, w2, b2)  
    a3 = dense(a2, w3, b3)  
    a4 = dense(a3, w4, b4)  
    return a4
```

$$\vec{w}_1^{[1]} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \vec{w}_2^{[1]} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}, \vec{w}_3^{[1]} = \begin{bmatrix} 5 \\ -6 \end{bmatrix}$$

$$\Rightarrow w = \text{np.array} \left( \begin{bmatrix} [1, -3, 5], \\ [2, 4, -6] \end{bmatrix} \right) \quad 2 \text{ by } 3$$

$$b_1^{[1]} = -1, b_2^{[1]} = 1, b_3^{[1]} = 2$$

$$\Rightarrow b = \text{np.array}([-1, 1, 2])$$

$$\vec{a}^{[0]} = \vec{x}$$

$$\bullet \Rightarrow a-in = \text{np.array}([-2, 4])$$

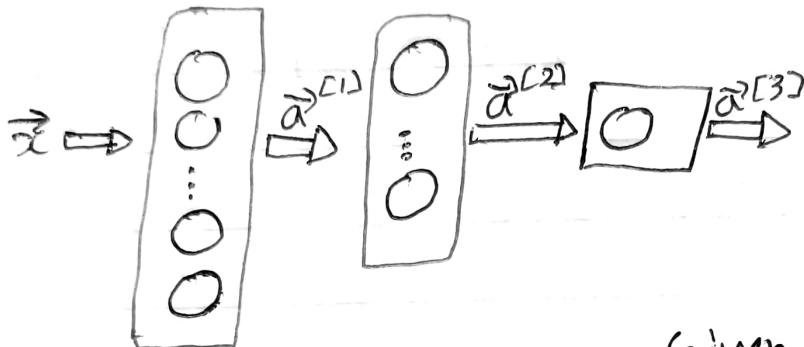
$\Rightarrow$  sequential(a-in) ↵

Now we know how to implement forward prop from scratch.  
☺

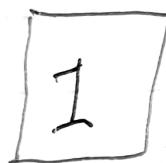
# Neural Network Training

week=2

## (i) Train a Neural Network in Tensorflow



Given set of  $(X, Y)$   
examples.



is 1 or not?

```
import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras import layers
```

```
model = keras.Sequential([  
    layers.Dense(units=25, activation="sigmoid"),  
    layers.Dense(units=15, activation="sigmoid"),  
    layers.Dense(units=1, activation="sigmoid")])
```

```
model.compile(loss=tf.keras.losses.BinaryCrossentropy())  
model.fit(x, y, epochs=100)
```

## (ii) Model Training Steps

- ① specify how to compute output given input  $x \in \mathbb{R}^4$  parameter  $w, b$  (define model)  
 $f_{\vec{w}, b}(\vec{x}) = ?$

logistic Regression

$$z = np.dot(x, w) + b$$

$$f_x = 1 / (1 + np.exp(-z))$$

Neural Network

model = Sequential([  
 Dense(...),  
 Dense(...),  
 Dense(...)])

- ② specify loss & cost

$$\text{Loss} = L(f_{\vec{w}, b}(\vec{x}), y)$$

$$\text{cost} = J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$$

$$\text{loss} = -y * np.log(f_x) - (1-y) * np.log(1-f_x)$$

model.compile(  
 -  
 -  
 -  
 - )

- ③ Train on data to minimize the cost function,  
 $J(\vec{w}, b)$

$$w = w - \alpha \nabla J(w)$$

$$b = b - \alpha \nabla J(b)$$

model.fit(  
 x, y, epochs=100  
 )

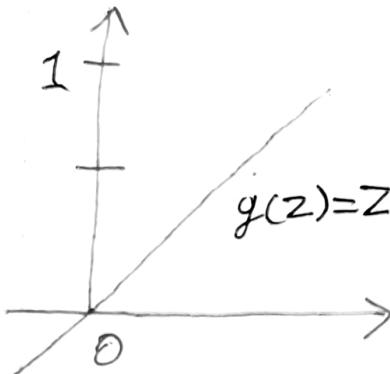
$$L = -y \log(f(\vec{x})) - (1-y) \log(1-f(\vec{x}))$$

same ↑ ↓  
 BinaryCrossEntropy() loss function for neural network

# Activation Functions

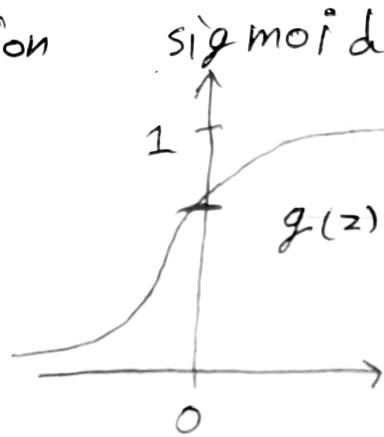
## (i) Alternatives to the sigmoid activation

Linear Activation Function



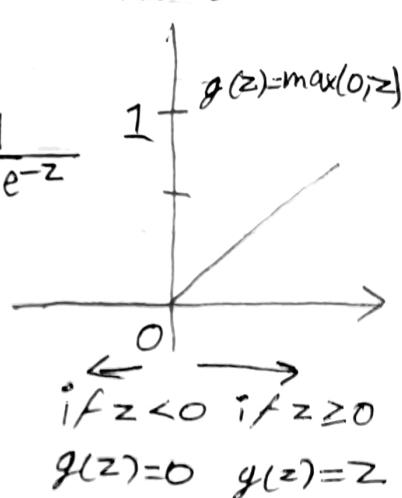
$$a = g(z) = \vec{w} \cdot \vec{x} + b$$

Sigmoid



$$0 < g(z) < 1$$

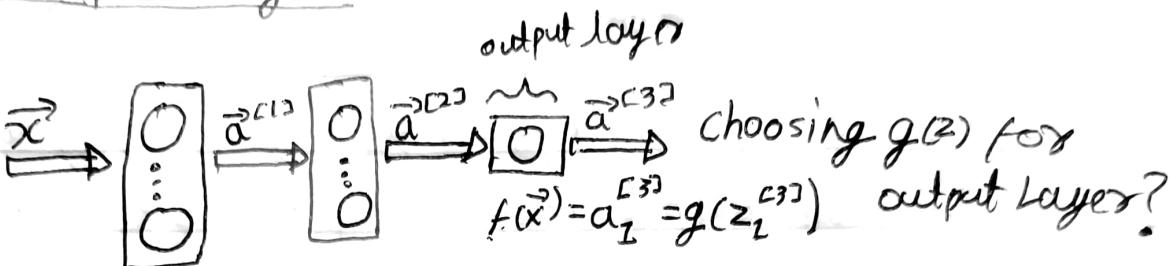
ReLU



Note:- ReLU activation function is widely used as it has a much lower run time, & runs much faster than any sigmoid function. but "ReLU function should only be used in the hidden layers."

## (ii) choosing activation functions

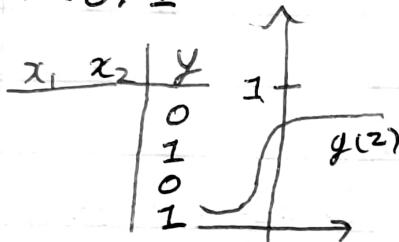
### Output Layer



Binary classification

⇒ sigmoid

$$Y = D/I$$

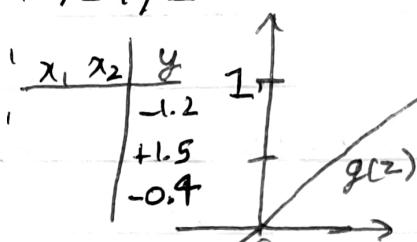


activation = "sigmoid"

Regression

Linear activation  $f^L$

$$Y = +/-$$

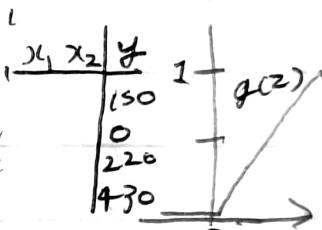


activation = "linear"

Regression

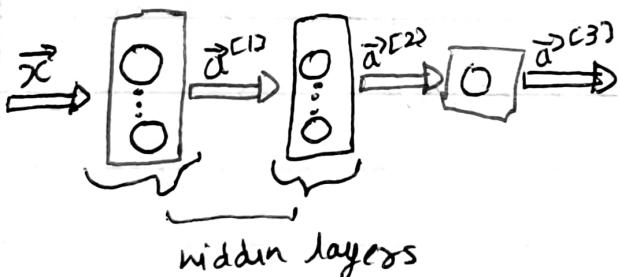
ReLU

$$Y = 0 \text{ or } +$$



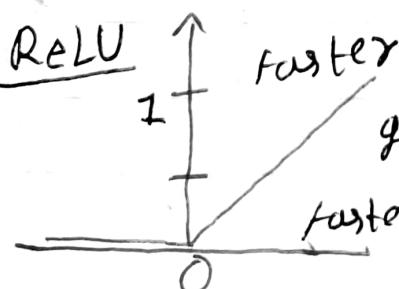
activation = "relu"

### Hidden Layer



choosing  $g(z)$  for hidden layers?

⇒ ReLU



more commonly  
choice rather than  
other activation  
function

# Multiclass Classification

## Softmax

Logistic Regression,  $y=0 \text{ or } 1$   
(2 possible output values)

$$\times a_1 = g(z) = \frac{e^{z_1}}{1+e^{-z_1}} = p(y=1|\vec{x})$$

$$O = 1 - a_1 = p(y=0|\vec{x})$$


---

softmax Regression,  $y=1, 2, \dots, N$   
( $N$  possible outputs)

$$z_j = \vec{w}_j \cdot \vec{x} + b_j, j=1, 2, \dots, N$$

parameters  $w_1, w_2, \dots, w_N$   
 $b_1, b_2, \dots, b_N$

$$a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = p(y=j|\vec{x})$$

note:  $a_1 + a_2 + a_3 + \dots + a_N = 1$

Softmax Regression,  $y=1, 2, 3, 4$   
(4 possible output values)

$$\begin{aligned} a_1 &= \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \\ &\quad \times O \quad \square \quad \Delta \\ &= p(y=1|\vec{x}) = 0.30 \end{aligned}$$

$$O z_2 = \vec{w}_2 \cdot \vec{x} + b_2,$$

$$\begin{aligned} a_2 &= \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \\ &= p(y=2|\vec{x}) = 0.20 \end{aligned}$$

$$\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3,$$

$$\begin{aligned} a_3 &= \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \\ &= p(y=3|\vec{x}) = 0.15 \end{aligned}$$

$$\Delta z_4 = \vec{w}_4 \cdot \vec{x} + b_4,$$

$$\begin{aligned} a_4 &= \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} \\ &= p(y=4|\vec{x}) = 0.35 \end{aligned}$$

## Cost

### Logistic regression

$$z = \vec{w} \cdot \vec{x} + b$$

$$a_1 = g(z) = \frac{1}{1 + e^{-z}} = P(y=1|\vec{x})$$

$$a_2 = 1 - a_1 = P(y=0|\vec{x})$$

$$\text{loss} = -y \log a_1 - (1-y) \log(1-a_1)$$

if  $y=1$                     if  $y=0$

$$J(\vec{w}, b) = \text{average loss}$$

### softmax regression

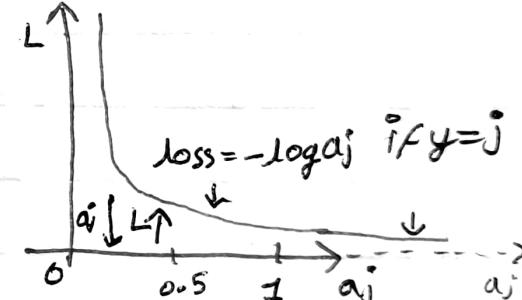
$$a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=1|\vec{x})$$

$$a_N = \frac{e^{z_N}}{e^{z_1} + e^{z_2} + \dots + e^{z_N}} = P(y=N|\vec{x})$$

### Crossentropy Loss

$$\text{loss}(a_1, \dots, a_N, y) =$$

$$\begin{cases} -\log a_1, & \text{if } y=1 \\ -\log a_2, & \text{if } y=2 \\ \vdots \\ -\log a_N, & \text{if } y=N \end{cases}$$



## Decision Trees

week = 3

## (i) Decision Tree Model

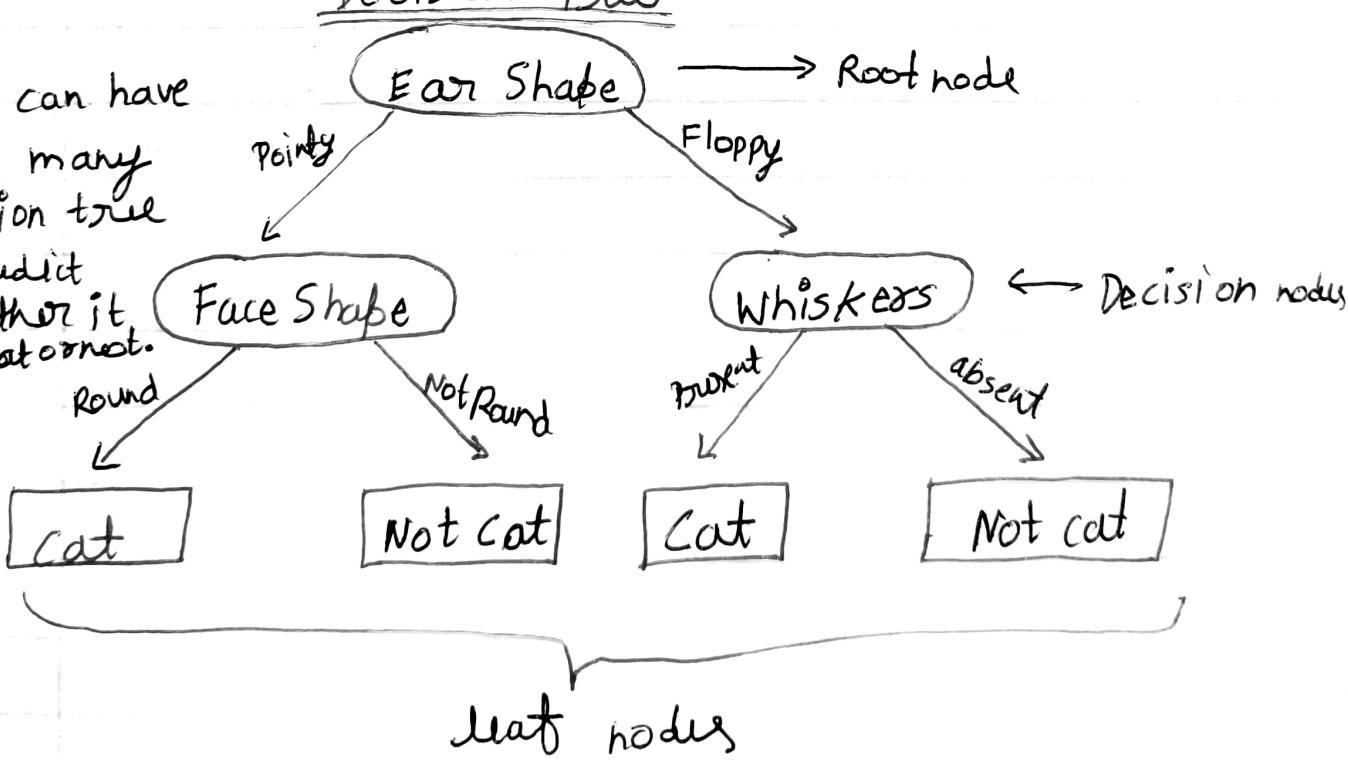
## Cat classification example

<sup>o</sup> images (x <sub>1</sub> )	<u>Ear shape</u>	<u>Face Shape (x<sub>2</sub>)</u>	<u>whiskers (x<sub>3</sub>)</u>	<u>cat</u>
I	Pointy	Round	Present	1
2	Floppy	not round	Present	1
3	Floppy	Round	Absent	0
:	:	:	:	:

categorical (discrete value).

## Decision Tree

we can have  
other many  
decision tree  
to predict  
whether it  
is a cat or not.

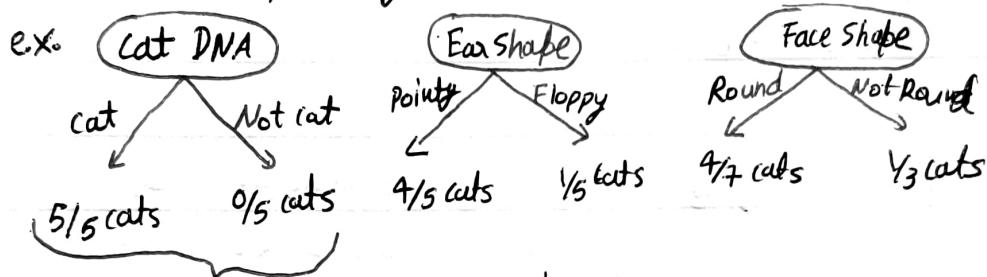


### (ii) Learning Process

#### Decision Tree Learning

Decision 1: How to choose what feature on at each node?

⇒ Maximize purity (or minimize impurity)



this is the maximize purity  
∴ we choose Cat DNA attribute

Decision 2: When do you stop splitting?

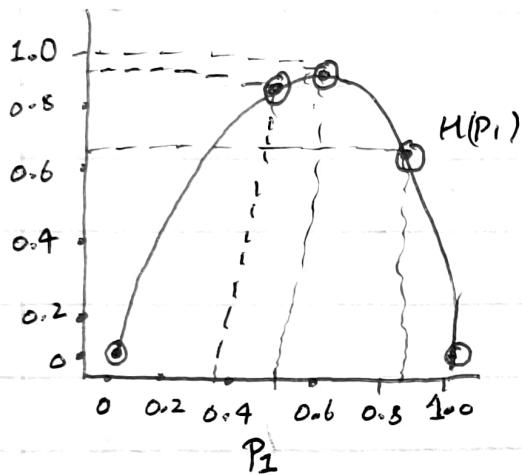
- When a node is 100% one class.
- When splitting a node will results in the exceeding a maximum depth.
- When improvements in purity score are below a threshold.
- When number of examples in a node is below a threshold.

# Decision Tree Learning

## (i) Measuring Purity

Entropy as a measure of impurity

$P_1$  = fraction of examples that are cats



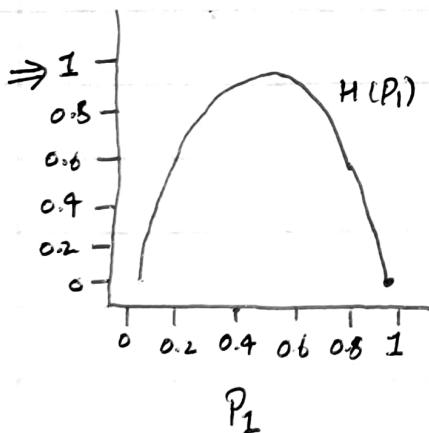
Example 1:  $P_1=0, H(P_1)=0.0$

Example 2:  $P_1=2/6, H(P_1)=0.42$

Example 3:  $P_1=3/6, H(P_1)=1.0$

Example 4:  $P_1=5/6, H(P_1)=0.65$

Example 5:  $P_1=6/6, H(P_1)=0.0$



$$\Rightarrow P_0 = 1 - P_1$$

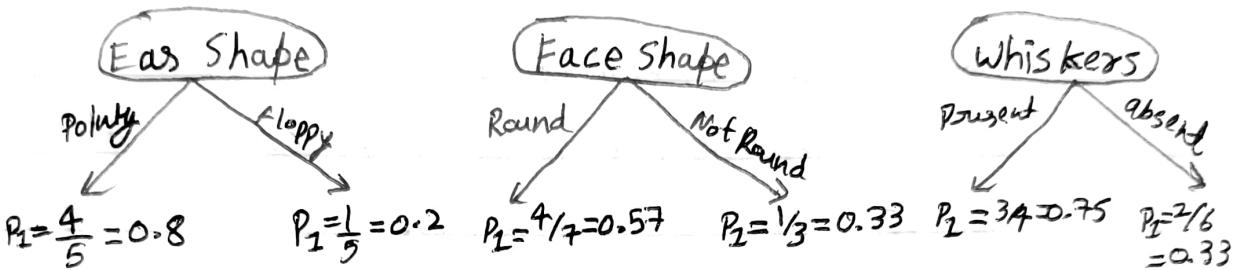
$$\Rightarrow H(P_1) = -P_1 \log_2(P_1) - P_0 \log_2(P_0)$$

$$= -P_1 \log_2(P_1) - (1-P_1) \log_2(1-P_1)$$

Note: " $0 \log(0) = 0$ "

## (ii) Choosing a split: Information Gain

### Choosing a split



$$H(0.5) = 1, \quad H(0.8) = 0.72 \quad H(0.2) = 0.72 \quad H(0.57) = 0.99 \quad H(0.33) = 0.92 \quad H(0.75) = 0.31 \quad H(0.33) = 0.92$$

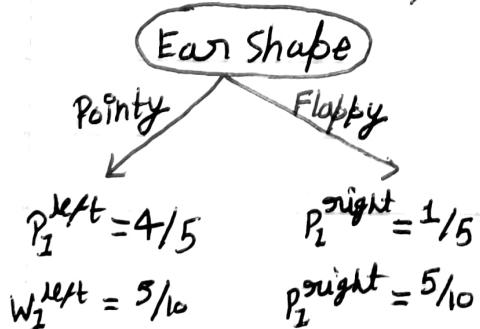
$$\Rightarrow H(0.5) - \left( \frac{5}{10} H(0.8) + \frac{5}{10} H(0.2) \right) \quad H(0.5) - \left( \frac{7}{10} H(0.57) + \frac{3}{10} H(0.33) \right) \quad H(0.5) - \left( \frac{4}{10} H(0.75) + \frac{6}{10} H(0.33) \right)$$

$$= 0.25 \text{ (Reduction Rate)} \quad = 0.03 \quad = 0.12$$

Information Gain

we choose this attribute  
as it has large entropy  
reduction rate.

Information Gain



Information Gain

$$= H(P_1^{\text{root}}) - \left( W^{\text{left}} H(P_1^{\text{left}}) + W^{\text{right}} H(P_1^{\text{right}}) \right)$$

### (iii) Putting it together

#### Decision Tree Learning

- start with all examples at the root node.
- calculate "information gain" for all possible features, & pick the one with the highest information gain.
- Split dataset according to selected feature, & create left & right branches of the tree.
- keep repeating splitting process until ~~stop~~ stopping criteria is met:
  - when a node is 100% one class.
  - when splitting a node will result in the tree exceeding a maximum depth.
  - Information gain from additional splits is less than threshold.
  - when number of examples in a node is below a threshold.

⇒ So, to do this we uses recursive splitting technique by which we first create left sub tree & right sub tree recursively.