

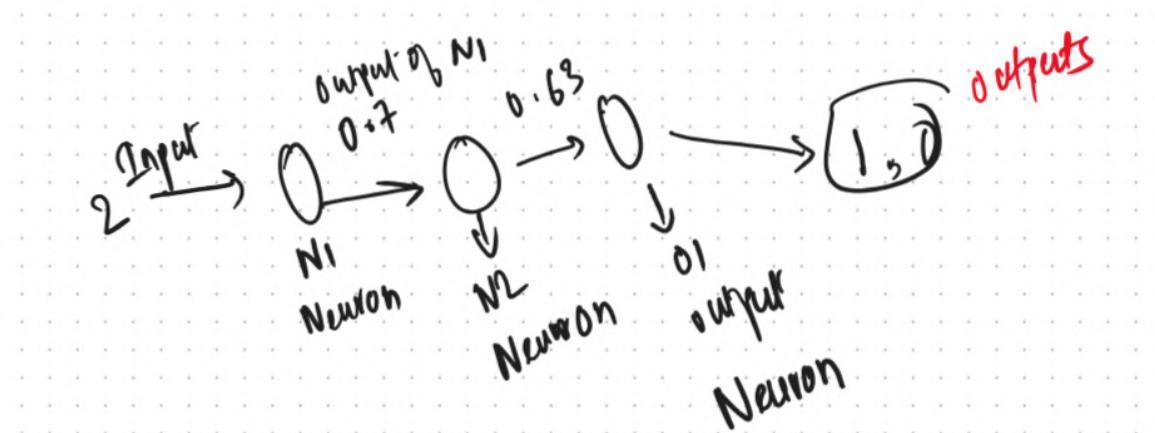
Advanced Learning Algorithms

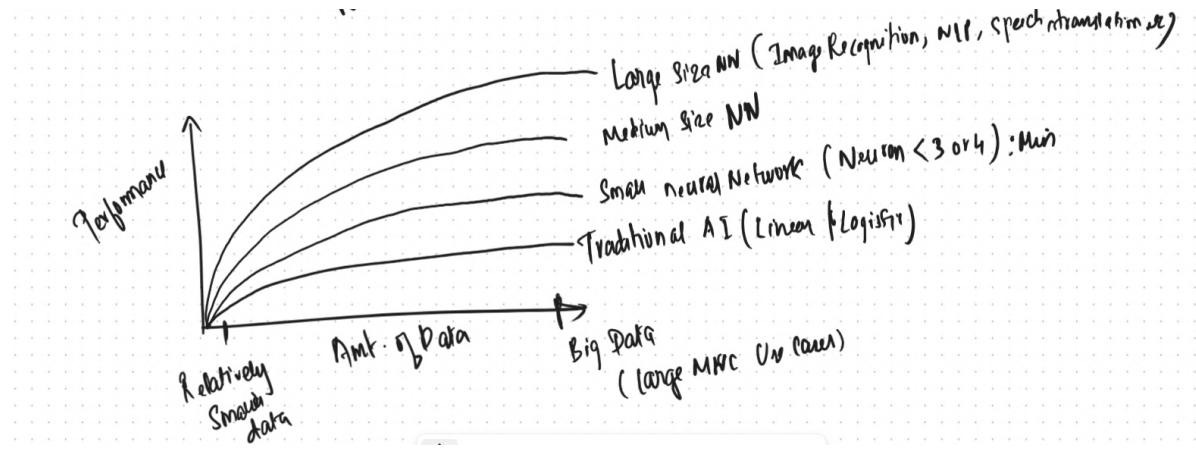
by Andrew NG

Author	Andrew NG
Link	https://www.coursera.org/learn/advanced-learning-algorithms/lecture/3wO3h/welcome
Score	
Status	In progress

▼ Neural Network Intuition

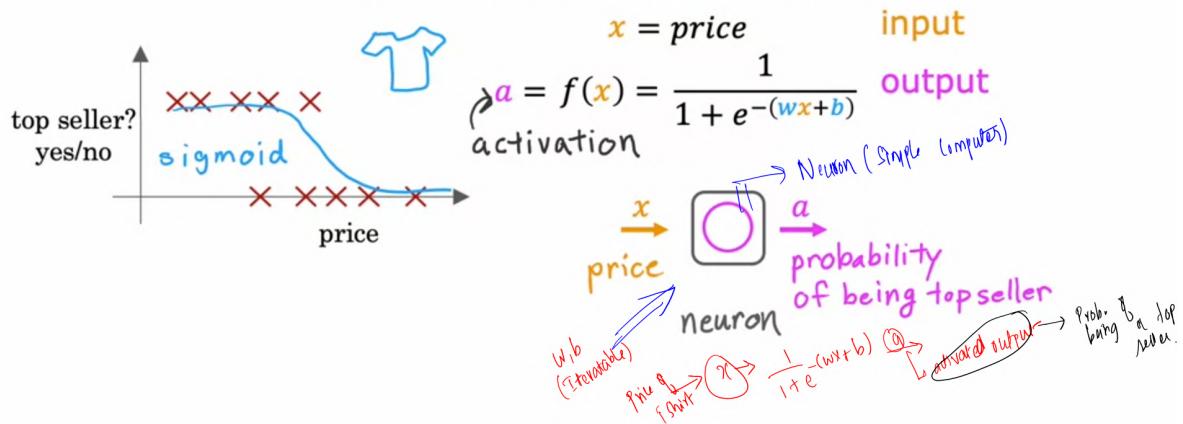
- Neural network was used to mimic the human brain by utilizing the software.
- the inputs are sent to neuron within which some calculations happens(functions and weights updates) post which we receive an output. This output then becomes an input for another neuron. Recent times we have been implying multiple neurons for our application. Though the neurons represent the surface level activities of the brain, it quiet effective for modern day application
-



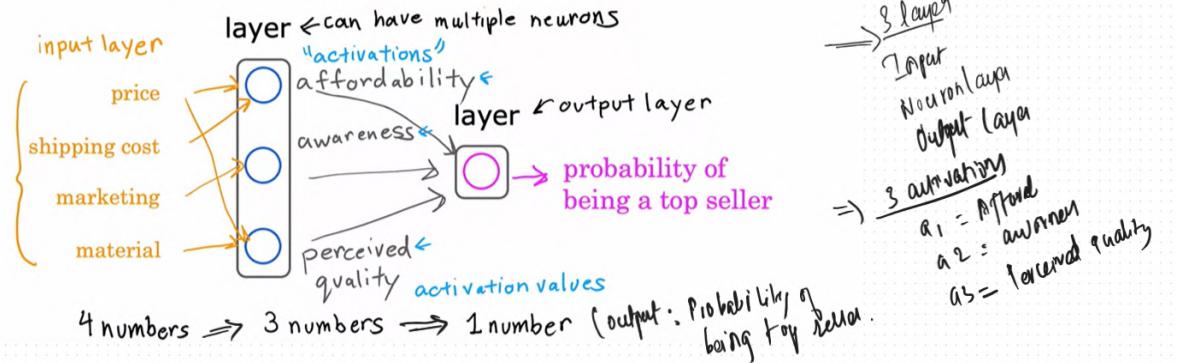


▼ Demand Prediction

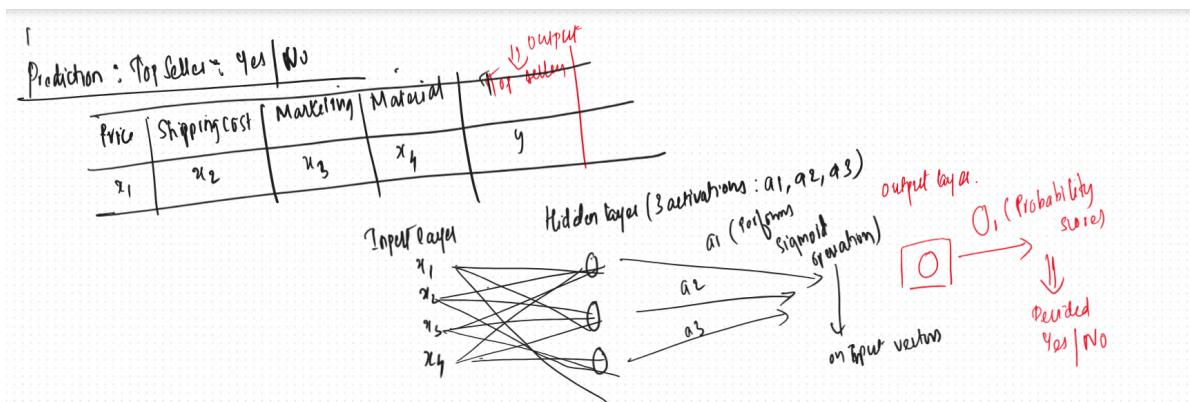
Demand Prediction



Demand Prediction

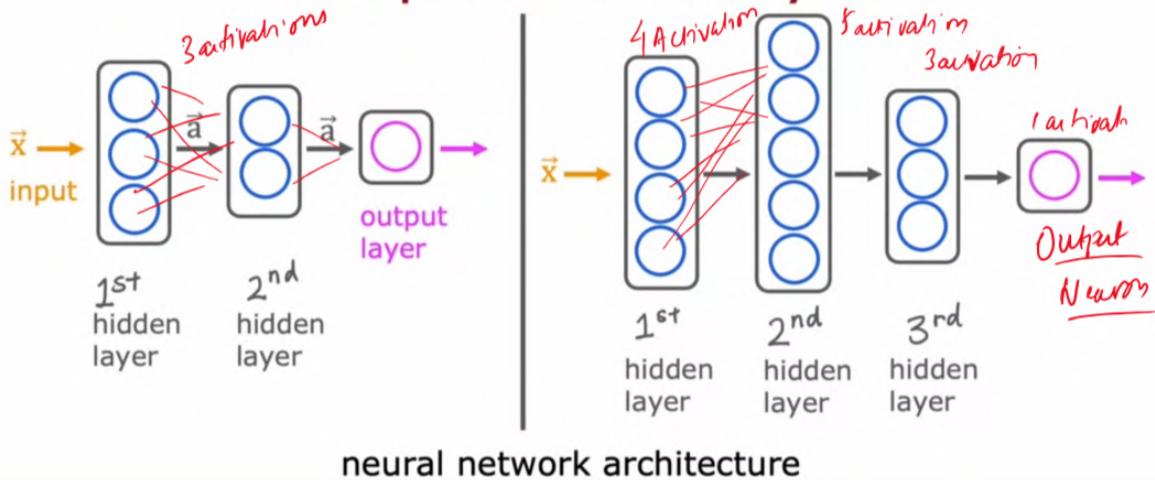


- In above diagram, we don't necessarily choose what input should be given to the neural layer as it becomes cumbersome process if there are many variables (say 100 variables), So we will connect all the input variables to all the neurons with neural layer.
- Each layer takes input vectors and gives output vectors
- Neural layers are called "**Hidden layers**".
- Earlier it was manual feature engineering, however here the neural layer performs feature engineering. Audacious thing here is that the output is not directly derived from input like linear or logistic regression, it is derived from hidden layer of 3 neurons.



- No feature engineering is required as hidden layer takes care of it. Manual feature selection might not ensure high quality of inputs to model, however hidden layer takes care of it.
- While handling the neural network, we should be able to decide "how many layers of neural network is needed" and "how many neurons per layer"
 - Answer is : neural network architecture

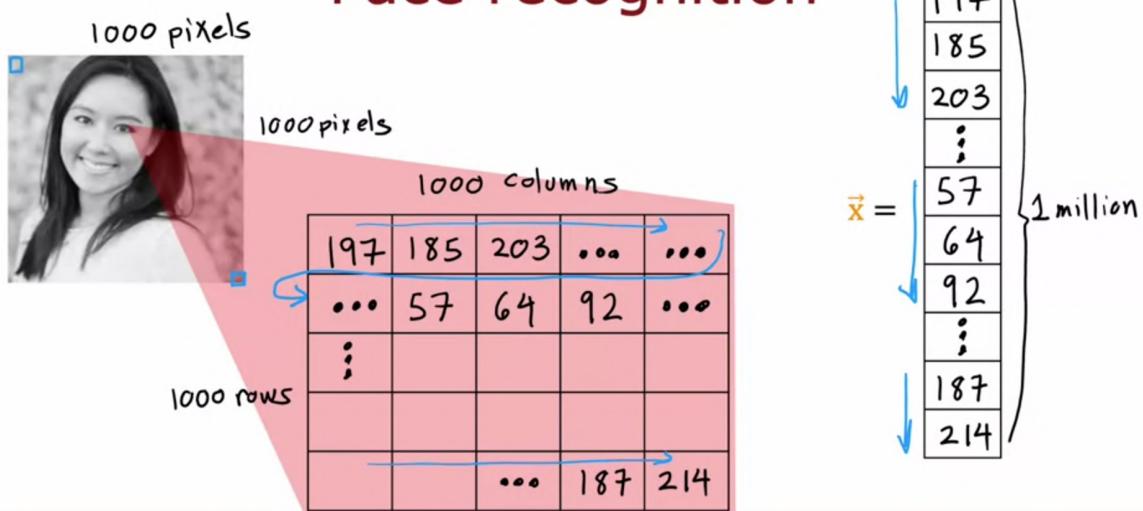
Multiple hidden layers



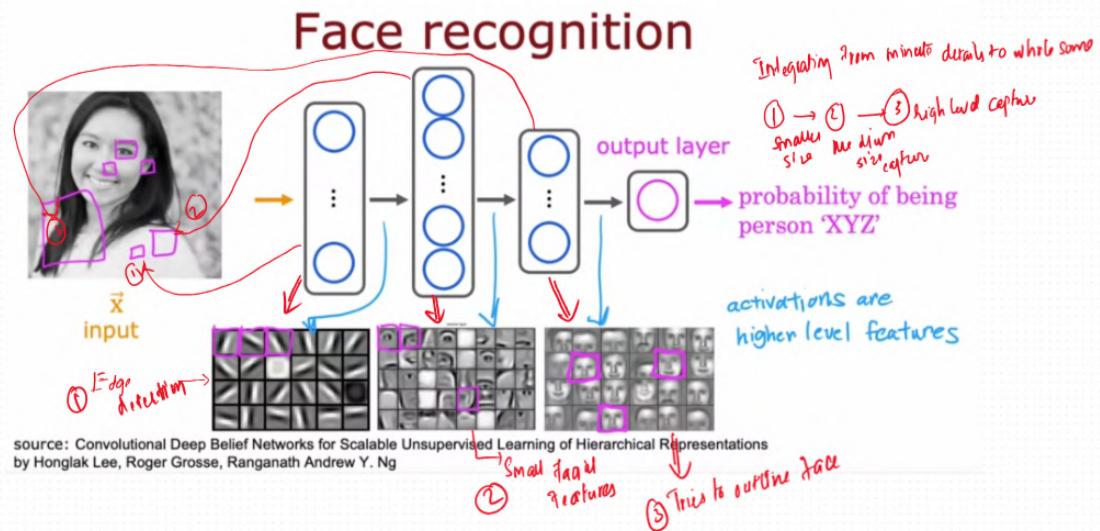
▼ Example : Face Recognition

- the images are recorded in the form of matrix. We know the image size. So Length * breadth of image = Number of pixels
 - Below 1000 rows * 1000 column pixel = 1M pixels
 - the intensity of light are captured, within this 1M pixels then they are unpacked to vectors in an array as stated below, which becomes input for hidden layers.

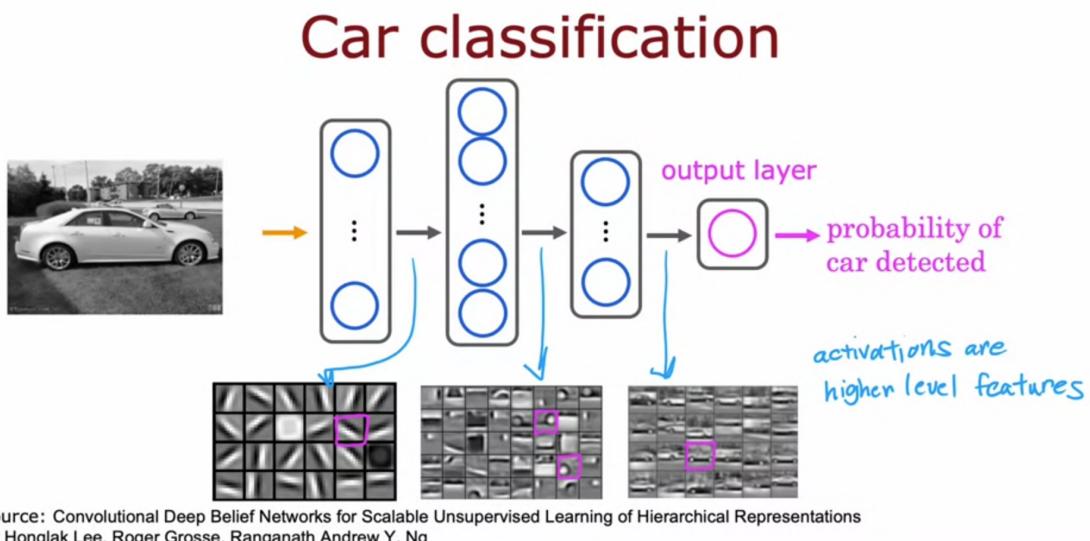
Face recognition



- In the below image, the first hidden layer is being visualized here we can find it tries to identify the edge of the image first.
- The beauty is that no one said or programmed to detect the edge first then smaller parts then whole face outline. It is done on its own from data.

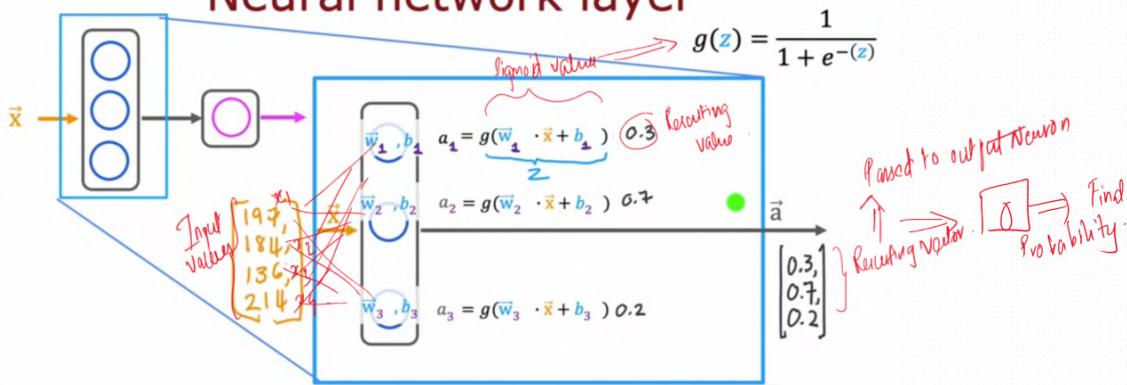


- Simply, changing the input data alone, the hidden layer learns on its own , the edges, smaller parts and finally a car. Here we are adding no extra hidden layer or extra neurons. The same model is applied to car as well.



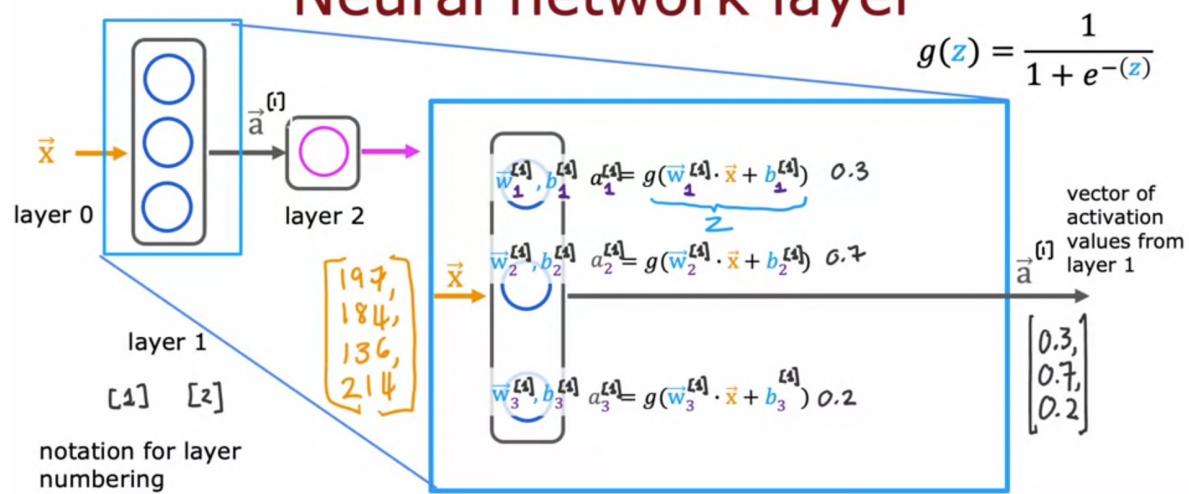
▼ Simple Neural Network Layer

Neural network layer

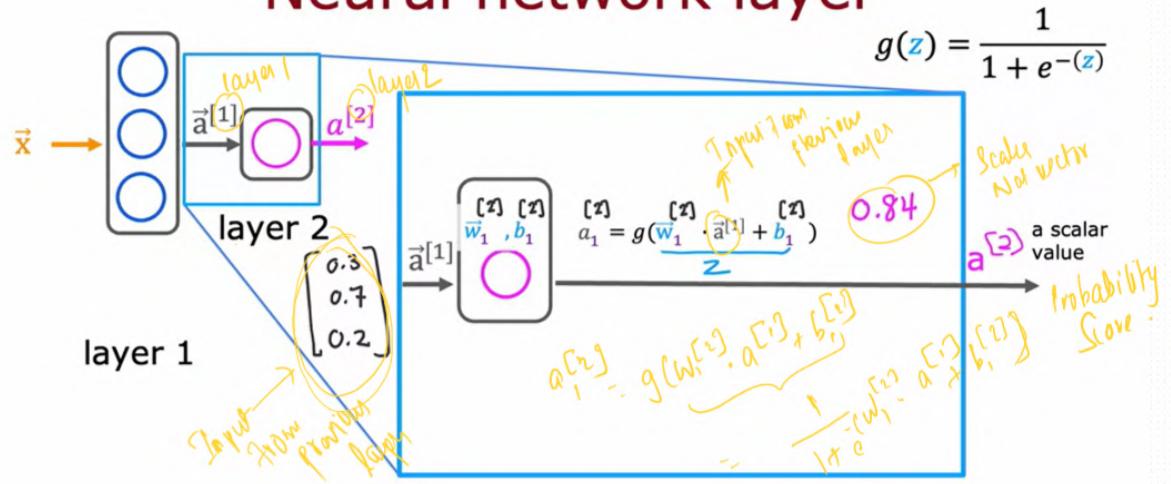


here the [1] denotes the nth hidden layer

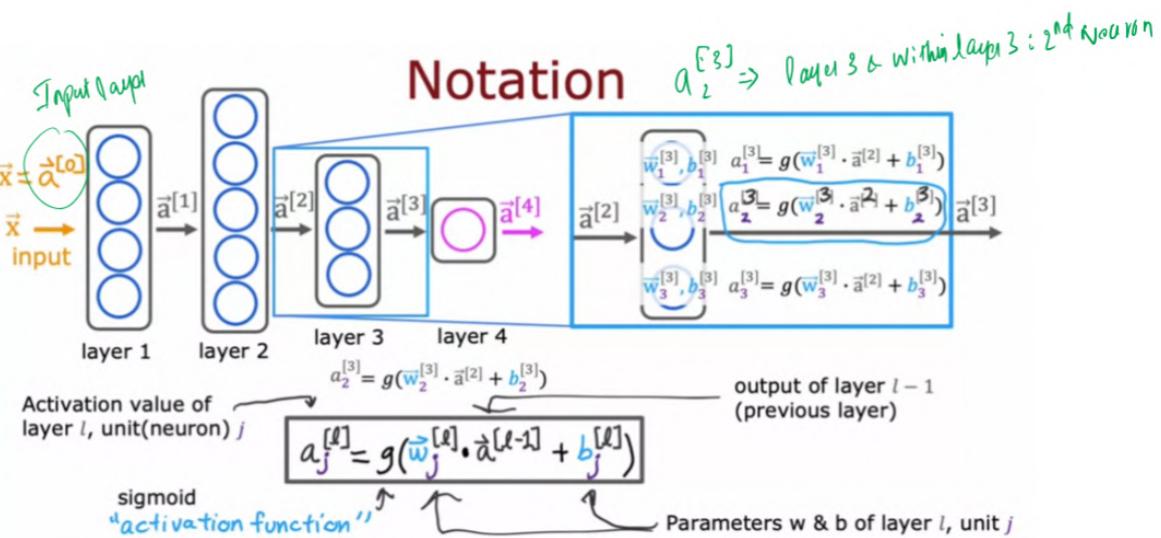
Neural network layer



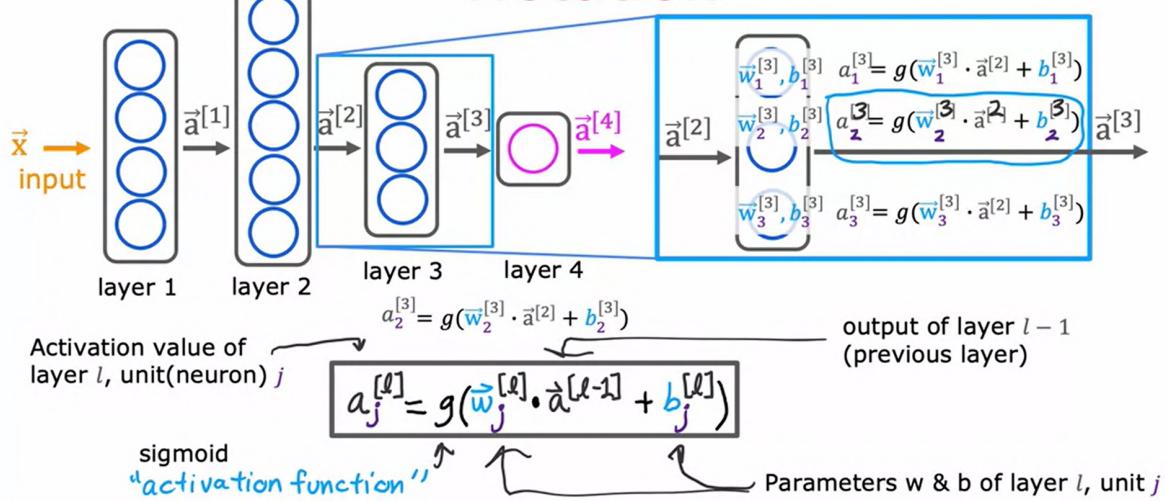
Neural network layer



▼ Complex Neural Network layer



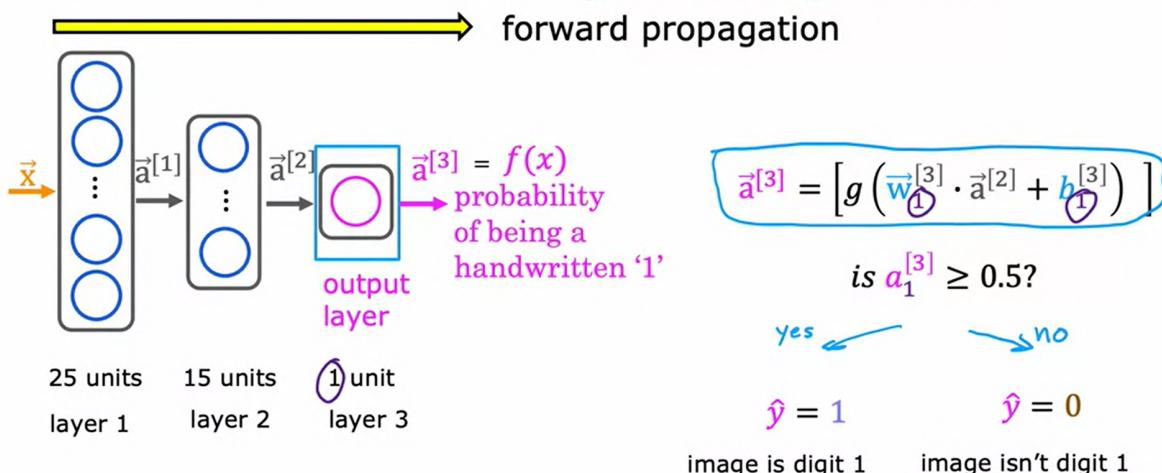
Notation



▼ Foreword propagation

- Here the input layer is taken followed by layer1 then layer 2 final output layer : which results in scalar, then we may use threshold to perform classification or we can use that prediction for regression. This sequential activity from left to right as stated below is called forward propagation. We can download the weights and neural network used for different application from the internet and use in our application.

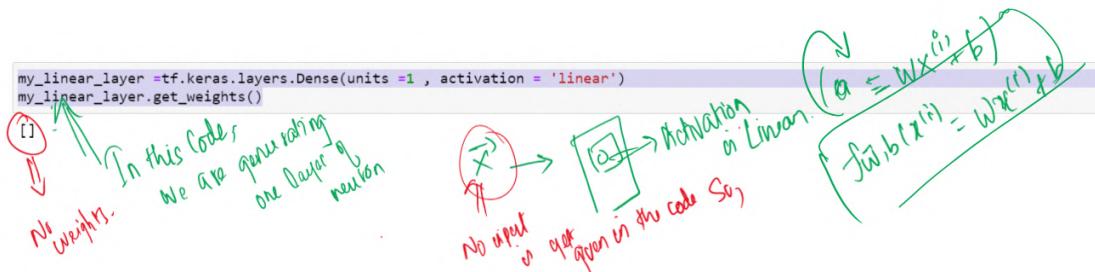
Handwritten digit recognition



- This is quiet fascinating that one model generated by larger MNC for a specific purpose can be used for different application by allowing it train

on different dataset.

- Tensorflow is a framework developed by google , keras(framework, a layer centric interface) is developed by an individual which later combined with tensorflow.
- A Neuron without activation(sigmoid, reLu etc.) is just “a Linear/ Regression Model”**
- A tensor is another name of array where the shape is in format of (1,1)**
- The initializing of weights is random.



- Input to linear layer

```
a1 = linear_layer(X_train[0].reshape(1,1))
print(a1)
```

```
tf.Tensor([[-1.03]], shape=(1, 1), dtype=float32)
```

- Get weights and set weights.

```
w, b= linear_layer.get_weights()
print(f"w = {w}, b={b}")
```

```
w = [[-1.03]], b=[0.]
```

- For setting the weights in linear_layer.

```

set_w = np.array([[200]])
set_b = np.array([100])

# set_weights takes a list of numpy arrays
linear_layer.set_weights([set_w, set_b])
print(linear_layer.get_weights())

[array([[200.]], dtype=float32), array([100.], dtype=float32)]

```

- 1000 pixels * 1000 pixels = 1 Million pixels this pixels is reshape to 2 dimensional array by applying **.reshape(-1,1)** which means the output will be **array[1Million,1]**.

```

print(np.array([0., 1, 2, 3, 4, 5], dtype=np.float32))
print(np.array([0., 1, 2, 3, 4, 5], dtype=np.float32).reshape(-1,1))

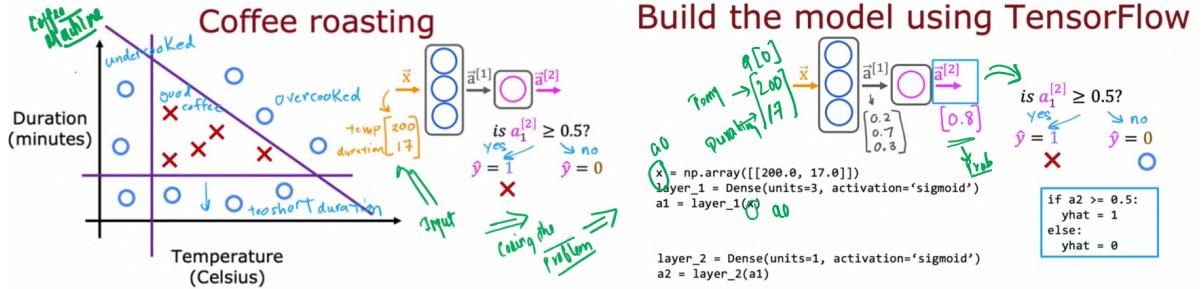
[0. 1. 2. 3. 4. 5.]
[[0.]
 [1.]
 [2.]
 [3.]
 [4.]
 [5.]]

```

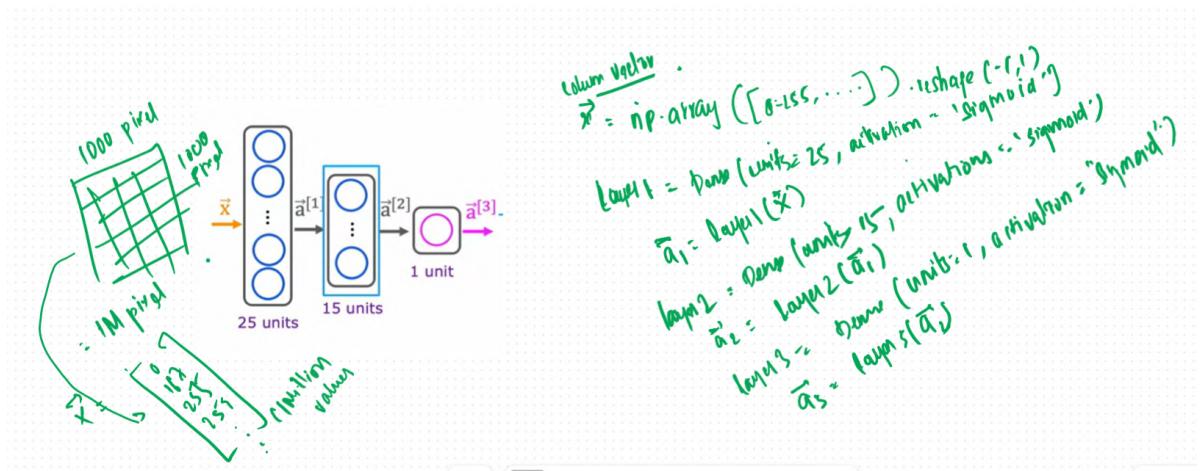
- Keras are generally very useful in creating multilayer models , the sequential function helps in achieving the same.

▼ How to implement tensor flow

- Tensorflow is tool to implement the deep learning projects.
- Other popular tool is pytorch.
- **Dense** : A type of neural layer or hidden layer.



- Simple implementation of tensorflow.

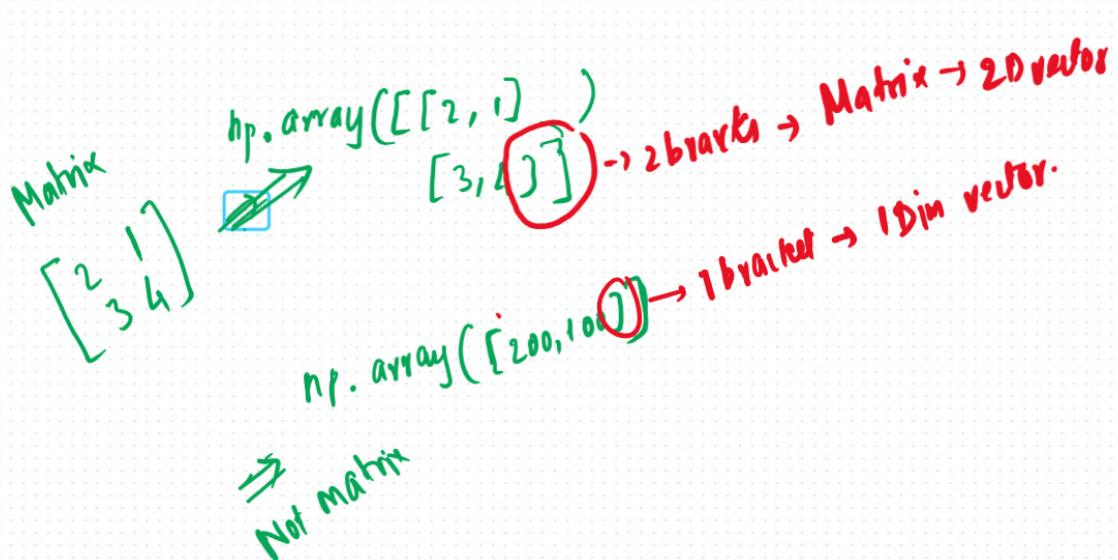


▼ How TensorFlow understands data

- Numpy was initial form of library which is standard for implementing the linear and logistic regressions
- Google brain created the tensor flow, so the contradiction could occur between numpy and tensorflow. Tensor is a datatype created by google brain team to carry out the function and data manipulation & handling. Tensor is a matrix a 2-Dim vector.
- Matrix → numpy array

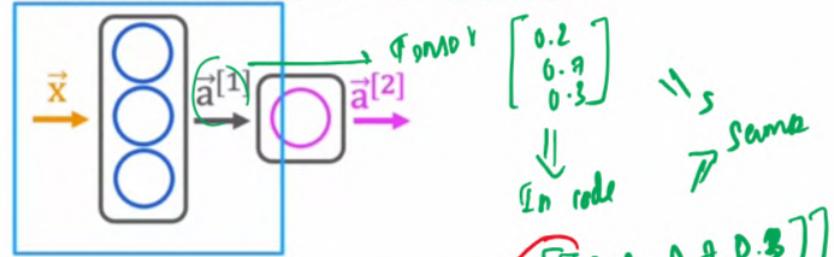
Note about numpy arrays

$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$	$x = np.array([[1, 2, 3], [4, 5, 6]])$	Matrix \rightarrow Array Conversion
2 rows \downarrow $\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}$	2 rows	2D representation
2 columns \downarrow $\begin{bmatrix} 1 & 2 \\ 4 & 5 \\ 6 & 7 \end{bmatrix}$	2D array	
2 rows \downarrow $\begin{bmatrix} 0.1 & 0.2 \\ -3 & -4 \\ -0.5 & -0.6 \\ 7 & 8 \end{bmatrix}$	2D array	2 x 3
4 rows \downarrow $\begin{bmatrix} 0.1 & 0.2 \\ -3 & -4 \\ -0.5 & -0.6 \\ 7 & 8 \end{bmatrix}$	$x = np.array([[0.1, 0.2], [-3, -4], [-0.5, -0.6], [7, 8]])$	4 x 2
2 columns \downarrow $\begin{bmatrix} 0.1 & 0.2 \\ -3 & -4 \\ -0.5 & -0.6 \\ 7 & 8 \end{bmatrix}$	$x = np.array([[0.1, 0.2], [-3, -4], [-0.5, -0.6], [7, 8]])$	1 x 2
4 x 2 matrix		2 x 1
	$\begin{matrix} R \times C \text{ in Matrix} \\ \text{Rows} \times \text{Columns} \end{matrix}$	$\begin{bmatrix} C_1 & C_2 & C_3 \\ R_1 & R_2 & R_3 \end{bmatrix}$



- Linear regression and logistic regression could handle the 1-Dim vector
- However for tensorflow the conventions are using the matrices which are 2-D vectors to represent the data. So this is computationally efficient in handling data.

Activation vector



```

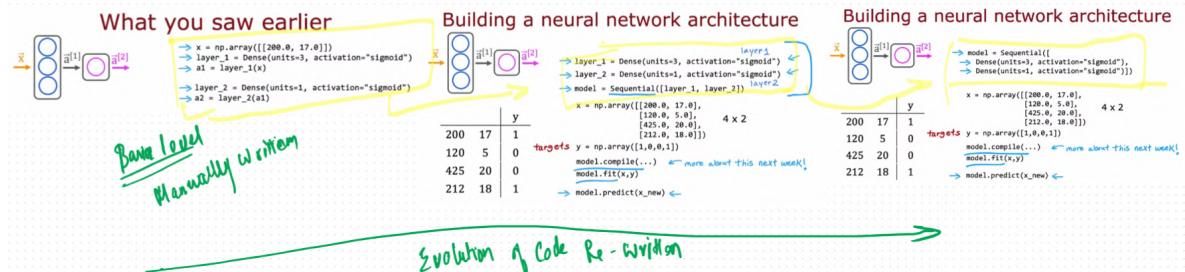
x = np.array([[200.0, 17.0]])
layer_1 = Dense(units=3, activation='sigmoid')
a1 = layer_1(x)
→ [[0.2, 0.7, 0.3]] 1 x 3 matrix
→ tf.Tensor([[0.2 0.7 0.3]], shape=(1, 3), dtype=float32)
→ a1.numpy() => to convert tensor → numpy array.
array([[0.2, 0.7, 0.3]], dtype=float32)

```

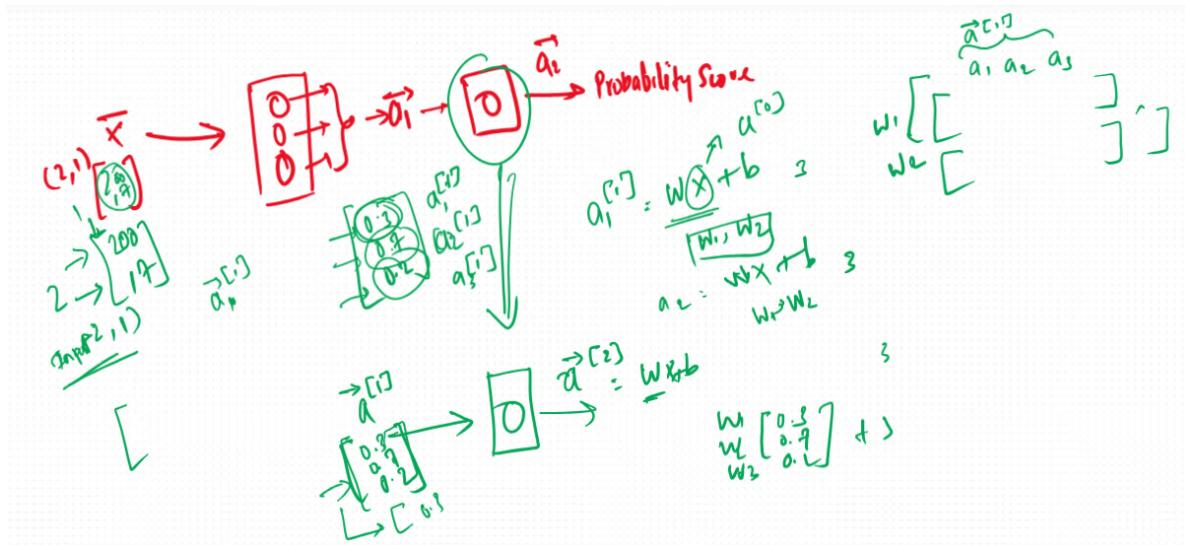
•

▼ How to build a TensorFlow neural network

- We are creation “**sequential**” function to avoid the carryout of activations to other layers manually, This will be taken care within the sequential function.



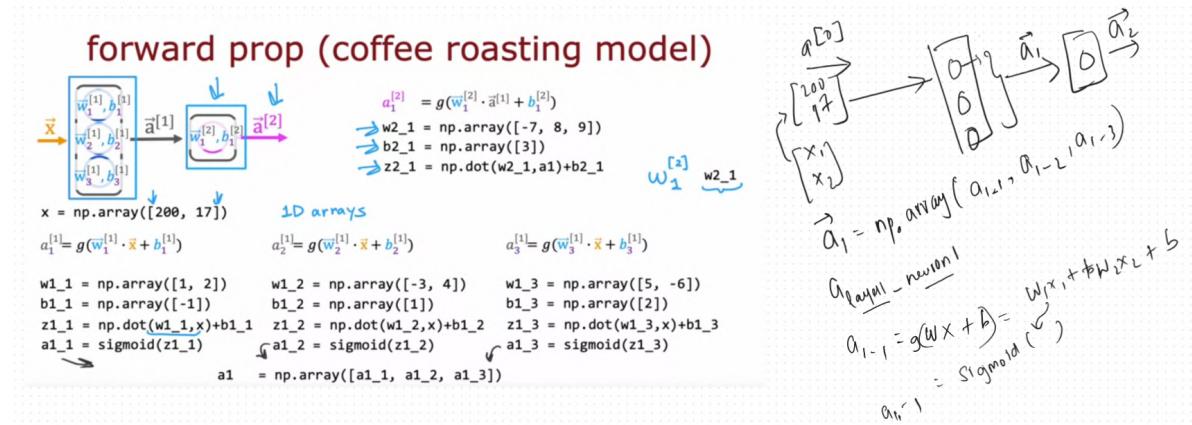
- The input shape imputed within model is not mandatory as this can be represented while fitting the model.
- Using Sigmoid function at the output layer is not good practice as**



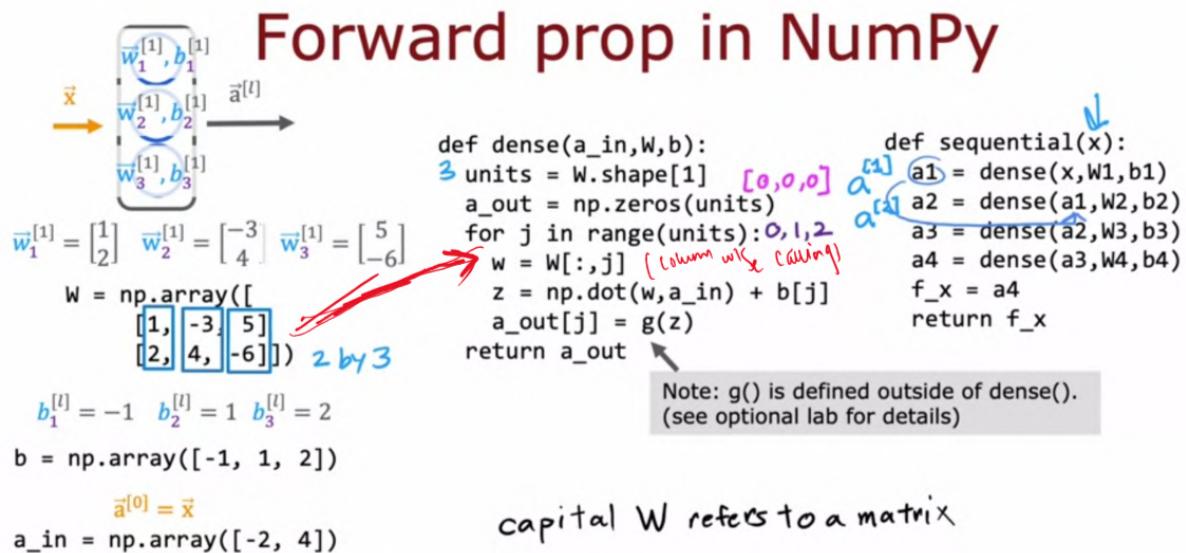
-
- The `model.compile` statement defines a loss function and specifies a compile optimization.
- The `model.fit` statement runs gradient descent and fits the weights to the data.
- What is epochs ?
 - Epochs defines the number of times the training set should be exposed to neural model.
 - if Epochs = 10, which means the whole training set is exposed to model 10 times.
 - these epochs are broken to batches, by default a batch in tensorflow consist of 32 units.
 - if 200,000 records, then batches would be $200,000/32 = 6250$ batches.
- Always each unit(neuron) in a layer tries to differentiate and identify. Different neurons within same layer has different methodology of identification. This is similar to average of marginal gain. **It is worth noting that the network learned these functions on its own through the process of gradient descent.** They are very much the same sort of functions a person might choose to make the same decisions.

Small neural network creation

▼ How to build the neural network forward propagation with python



▼ How to build the forward propagation simplistically



▼ What is AGI ?

- Artificial General Intelligence
- ANI : Artificial narrow intelligence : specific domain only

▼ What is vectorization, how is it efficient?

•

For loops vs. vectorization

```
x = np.array([200, 17])
W = np.array([[1, -3, 5],
              [-2, 4, -6]])
b = np.array([-1, 1, 2])

def dense(a_in,W,b):
    units = W.shape[1]
    a_out = np.zeros(units)
    for j in range(units):
        w = W[:,j]
        z = np.dot(w, a_in) + b[j]
        a_out[j] = g(z)
    return a_out
```

[[1,0,1]]

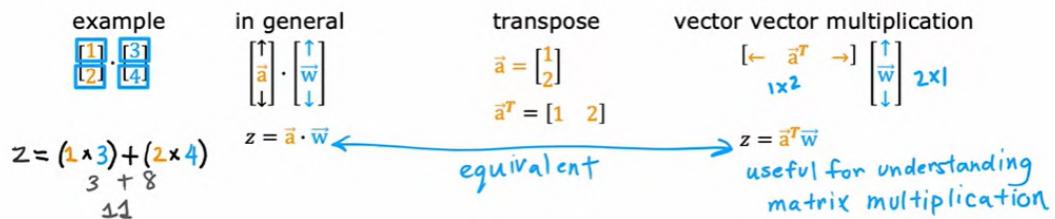
```
X = np.array([[200, 17]]) 2Darray
W = np.array([[1, -3, 5], same
              [-2, 4, -6]])
B = np.array([-1, 1, 2]) 1x3 2Darray
def dense(A_in,W,B): all 2Darrays
    Z = np.matmul(A_in,W) + B
    A_out = g(Z) matrix multiplication
    return A_out
```

↑ Matrix (computationally faster & efficient)

- How does this matmul works ?

- how dot product works ? : **Vector to vector multiplication**

Dot products



- how vector matrix multiplication occurs ? **vector to matrix multiplication**

Vector matrix multiplication

$$\vec{a} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\vec{a}^T = [1 \ 2] \quad W = \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix}$$

$$Z = \vec{a}^T W \quad [\leftarrow \vec{a}^T \rightarrow] \quad \begin{bmatrix} \uparrow & \uparrow \\ \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow \end{bmatrix}$$

1 by 2

$$Z = [\vec{a}^T \vec{w}_1 \quad \vec{a}^T \vec{w}_2]$$

$$(1 * 3) + (2 * 4) \quad (1 * 5) + (2 * 6)$$

$$3 + 8 \quad 5 + 12$$

$$11 \quad 17$$

$$Z = [11 \ 17]$$

- **What is a matrix ?** : When multiple vectors compiled together in columns to form matrix.
 - However when a matrix is transposed then the vectors are arranged in row wise.
- How matrix to matrix multiplication occurs ?

matrix matrix multiplication

$$A = \begin{bmatrix} 1 & -1 \\ 2 & -2 \end{bmatrix}$$

$$A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \end{bmatrix}$$

rows

$$W = \begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix}$$

columns

$$Z = A^T W = \left[\begin{array}{c|c} \leftarrow \vec{a}_1^T & \rightarrow \\ \leftarrow \vec{a}_2^T & \rightarrow \end{array} \right] \begin{bmatrix} \uparrow & \uparrow \\ \vec{w}_1 & \vec{w}_2 \\ \downarrow & \downarrow \end{bmatrix}$$

$$\begin{array}{c} \text{row1 col1} \\ \text{row2 col1} \end{array} = \begin{bmatrix} \vec{a}_1^T \vec{w}_1 & \vec{a}_1^T \vec{w}_2 \\ \vec{a}_2^T \vec{w}_1 & \vec{a}_2^T \vec{w}_2 \end{bmatrix} \begin{array}{c} \text{row1 col2} \\ \text{row2 col2} \end{array}$$

$$\begin{array}{c} (-1 \times 3) + (-2 \times 4) \\ -3 + -8 \\ -11 \end{array} \quad \begin{array}{c} (-1 \times 5) + (-2 \times 6) \\ -5 + -12 \\ -17 \end{array}$$

$$= \begin{bmatrix} 11 & 17 \\ -11 & -17 \end{bmatrix}$$

general rules for
matrix multiplication
↳ next video!

Matrix multiplication rules

$$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix} \quad Z = A^T W = \begin{bmatrix} 11 & 17 & 23 & 9 \\ -11 & -17 & -23 & -9 \\ 1.1 & 1.7 & 2.3 & 0.9 \end{bmatrix}$$

$$\vec{a}_1^T \vec{w}_1 = (1 \times 3) + (2 \times 4) = 11$$

3 by 4 matrix

row 3 column 2

$$\vec{a}_3^T \vec{w}_2 = (0.1 \times 5) + (0.2 \times 6) = 1.7$$

0.5 + 1.2

row 2 column 3?

$$\vec{a}_2^T \vec{w}_3 = (-1 \times 7) + (-2 \times 8) = -23$$

-7 + -16

Matrix multiplication rules

$$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix} \quad Z = A^T W = \begin{bmatrix} 11 & 17 & 23 & 9 \\ -11 & -17 & -23 & -9 \\ 1.1 & 1.7 & 2.3 & 0.9 \end{bmatrix}$$

3 x 2 2 x 4

3 by 4 matrix

can only take dot products
of vectors that are same length

$$\begin{bmatrix} 0.1 & 0.2 \end{bmatrix} \quad \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

length 2 length 2

- ▼ if vectorization is about matrix creation, then how to multiple using python codes ?

- Here **T** is the transpose function:

Matrix multiplication in NumPy

$$A = \begin{bmatrix} 1 & -1 & 0.1 \\ 2 & -2 & 0.2 \end{bmatrix} \quad A^T = \begin{bmatrix} 1 & 2 \\ -1 & -2 \\ 0.1 & 0.2 \end{bmatrix} \quad W = \begin{bmatrix} 3 & 5 & 7 & 9 \\ 4 & 6 & 8 & 0 \end{bmatrix} \quad Z = A^T W = \begin{bmatrix} 11 & 17 & 23 & 9 \\ -11 & -17 & -23 & -9 \\ 1.1 & 1.7 & 2.3 & 0.9 \end{bmatrix}$$

`A=np.array([[1,-1,0.1], [2,-2,0.2]])` `W=np.array([[3,5,7,9], [4,6,8,0]])` `Z = np.matmul(AT,W)` *or* `Z = AT @ W`

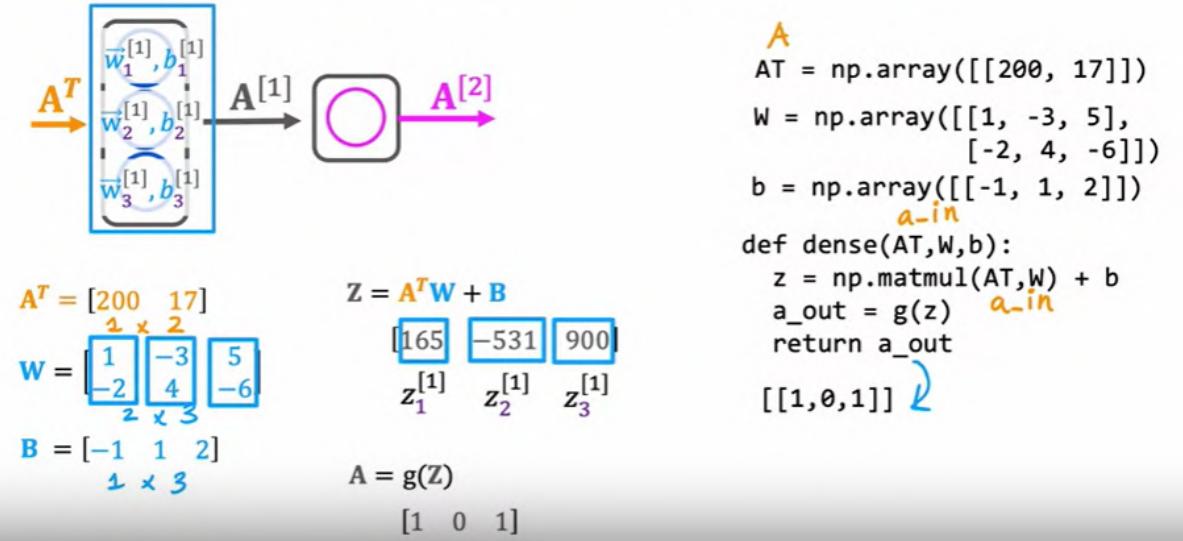
`AT=np.array([[1,2], [-1,-2], [0.1,0.2]])`

`AT=A.T` *transpose*

result $\begin{bmatrix} [11,17,23,9], [-11,-17,-23,-9], [1.1,1.7,2.3,0.9] \end{bmatrix}$

- Always remember when vectors are stacked in column wise that becomes the matrix, same happened here with W.

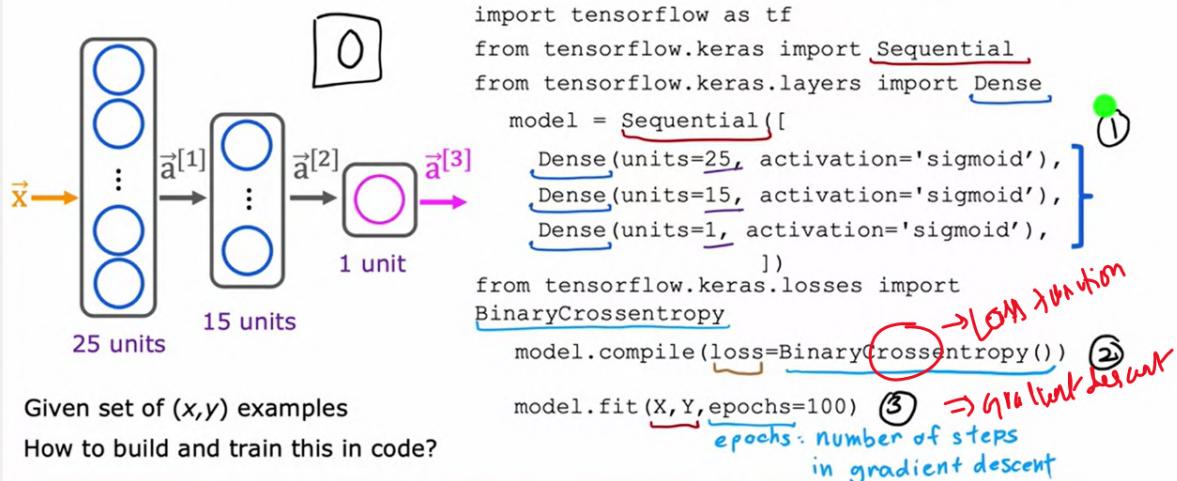
Dense layer vectorized



▼ How to train neural network ?

- First step : create a neural architecture using sequential function.
- Second step : model. Compile : explore the loss function
- 3rd Step : `model.fit` : apply gradient descent where the epochs is number of iteration taken within the gradient descent

Train a Neural Network in TensorFlow



- Analogies of logistic regression and neural network :

Model Training Steps Tensor Flow

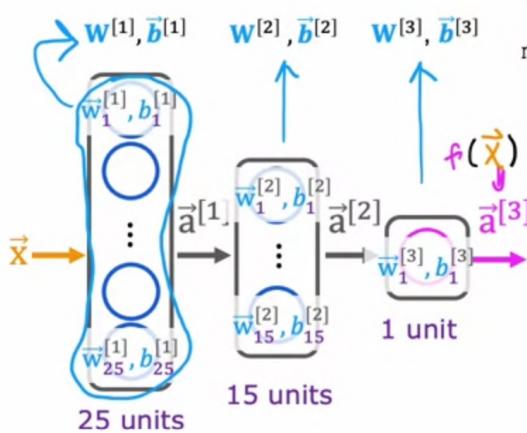
	logistic regression	neural network
① specify how to compute output given input x and parameters w, b (define model) $f_{\vec{w}, b}(\vec{x}) = ?$	$z = \text{np.dot}(w, x) + b$ $f_x = 1 / (1 + \text{np.exp}(-z))$	<code>model = Sequential([Dense(...), Dense(...), Dense(...)])</code>
② specify loss and cost $L(f_{\vec{w}, b}(\vec{x}), y)$ 1 example $J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$	logistic loss $\text{loss} = -y * \text{np.log}(f_x) - (1-y) * \text{np.log}(1-f_x)$ $w = w - \alpha * dj_dw$ $b = b - \alpha * dj_db$	binary cross entropy <code>model.compile(loss=BinaryCrossentropy())</code> <code>model.fit(X, y, epochs=100)</code>
③ Train on data to minimize $J(\vec{w}, b)$		

- Step 1 :

1. Create the model

define the model

$$f(\vec{x}) = ?$$



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

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='sigmoid'),
])
```

- Step 2 : Mentioning the loss function for the neuron , if classification and binary then we use binary cross entropy function, if they are regression problem then we might use mean squared error as loss function.

2. Loss and cost functions

handwritten digit classification problem $\xrightarrow{\text{binary classification}}$

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

Compare prediction vs. target

logistic loss

also Known as binary cross entropy

model.compile(loss= BinaryCrossentropy())
 regression
 (predicting numbers and not categories)
 model.compile(loss= MeanSquaredError())

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)})$$

$w^{[1]}, w^{[2]}, w^{[3]}$ $b^{[1]}, b^{[2]}, b^{[3]}$

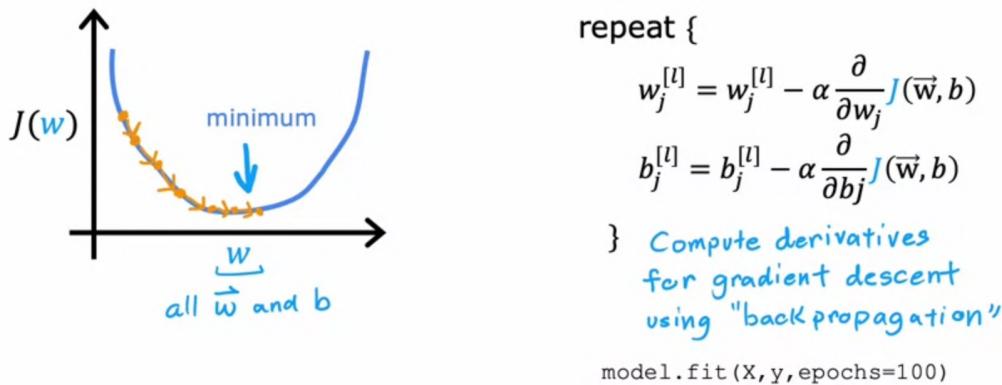
$f_{W,B}(\vec{x})$

from tensorflow.keras.losses import
 BinaryCrossentropy Keras

from tensorflow.keras.losses import
 MeanSquaredError

- Step 3 : fitting the model : this does not limit to just applying gradient descent : this gradient descent is achieved by back propagation > while passing backwards the weights are updated thus resulting in updation of weights causing the application of gradient descent.

3. Gradient descent

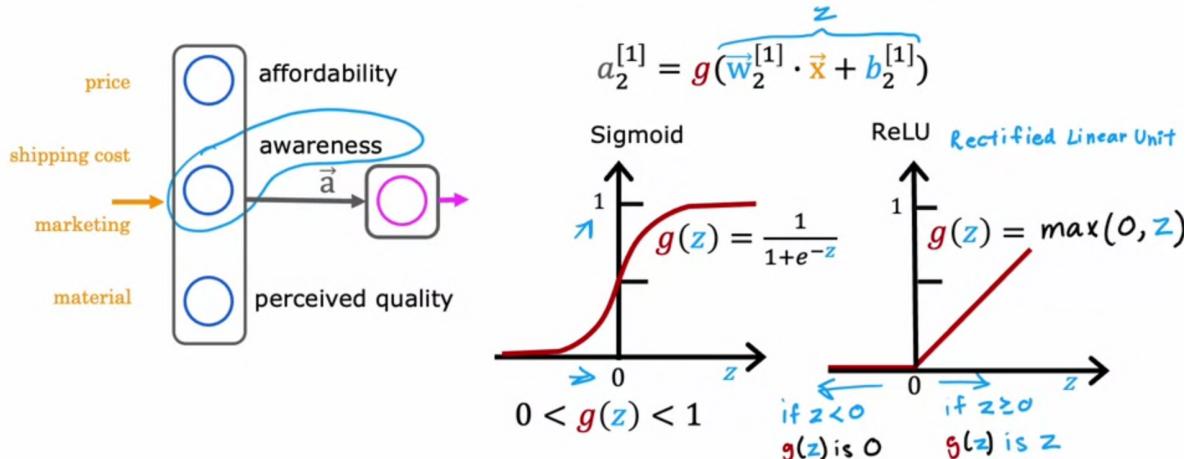


- Now a days, developers are using the libraries for implementing the deep learning rather than building it from scratch.

▼ What are different types of activation functions ?

- ReLU : Rectified linear unit : Just a naming conventions used within the industry:

Demand Prediction Example

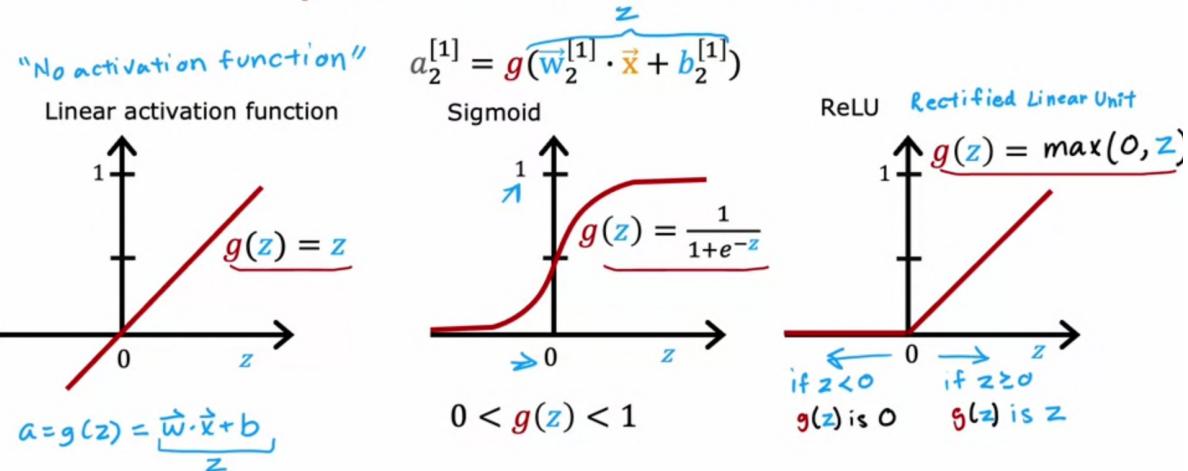


- In Linear activation : this is completely different from not using an activation function, however people sometime refer to this as no activation

function. But this is great point to start with.

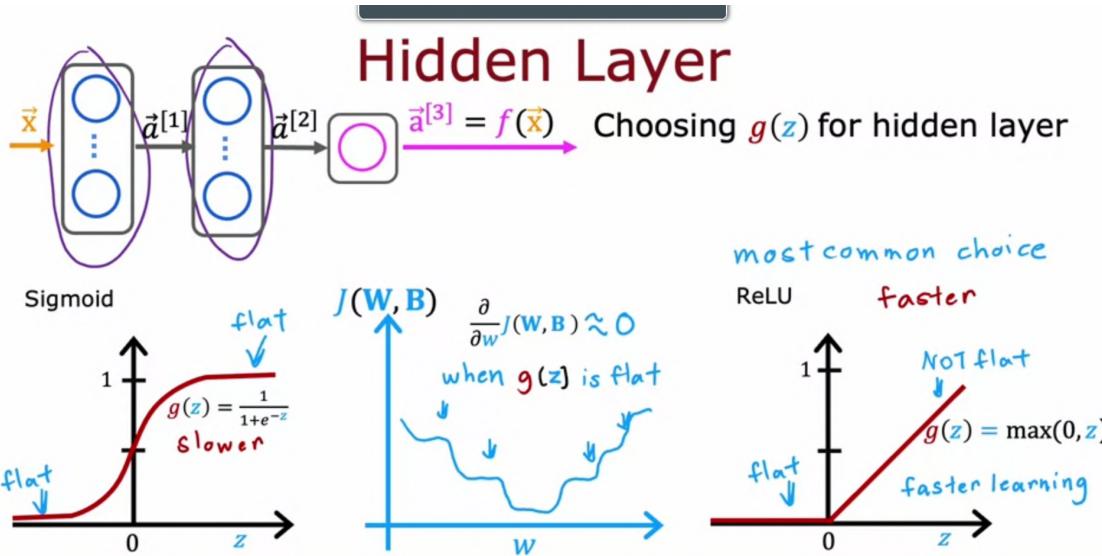
- Major activation function within the market :

Examples of Activation Functions



- ▼ How do we select the better activation functions for your application ?

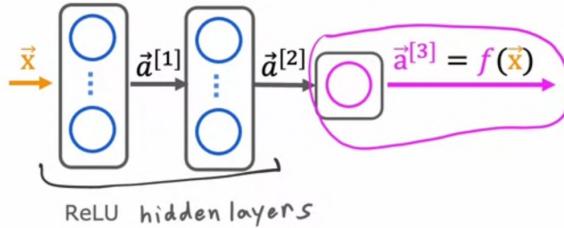
- For binary classification function : **Sigmoid** classification function
- For regression problem : **Linear activation** would fit this case (this can be applied if the output can contain negative values too eg. Stock prices etc.)
- For regression problem : with non negative output > like house price prediction : **ReLU** would better fit the problem.
- For multiclass classification : **Softmax** function
- For activation within the neural layer or hidden layer : popular choice here is **ReLU** rather than a sigmoid function. But why ?
 - **ReLU is computation is faster** as it has only $\max(0, z)$ however for sigmoid it has $1/(1+e^{-z})$
 - ReLU is flat on only one side meaning, the value does not change for only one side, however in sigmoid function both the sides are constrained with boundaries. So the gradient descent would be slower for sigmoid than the ReLU function.



- o **Recommendation**

- Output based
 - if classification : Use sigmoid function
 - if regression : use Linear based function and ReLU for non negative outputs
 - if multiclass classification : you can used softmax function
- Hidden layer
 - Use ReLU function for better performance.
-

Choosing Activation Summary



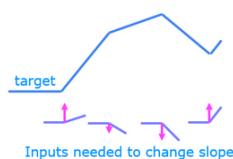
binary classification
activation='sigmoid'
regression $y \leq 0$
activation='linear'
regression $y \geq 0$
activation='relu'

```
from tensorflow import keras
model = Sequential([
    Dense(units=25, activation='relu'), layer1
    Dense(units=15, activation='relu'), layer2
    Dense(units=1, activation='sigmoid') layer3
])
or 'linear'
or 'relu'
```

Other Activation function:

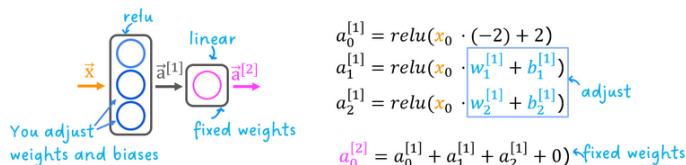
- Leaky ReLU activation functions, TanH activations functions , Switch activation functions

Why Non-Linear Activations?



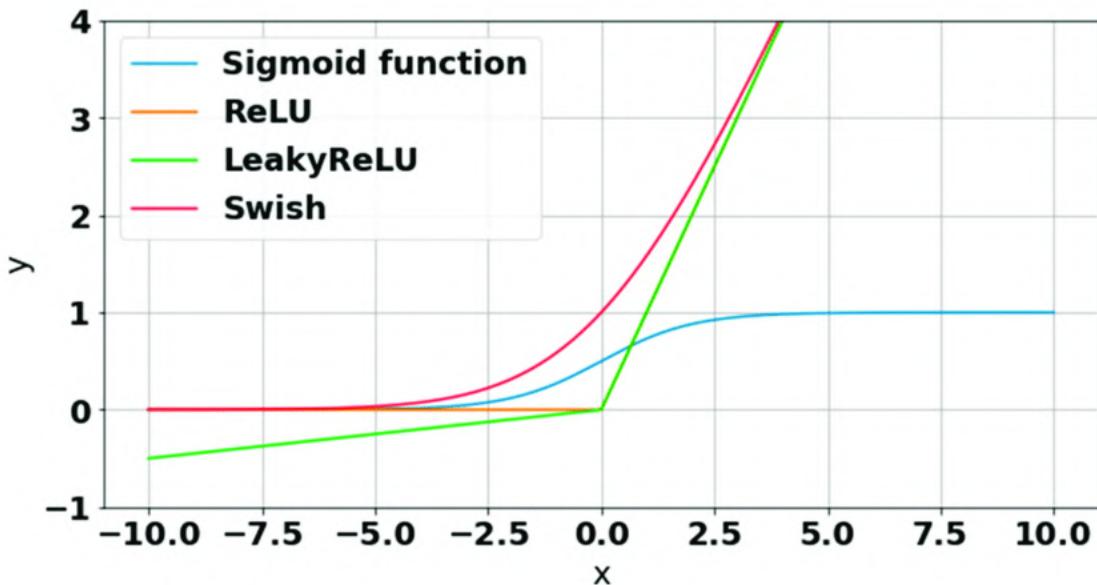
The function shown is composed of linear pieces (piecewise linear). The slope is consistent during the linear portion and then changes abruptly at transition points. At transition points, a new linear function is added which, when added to the existing function, will produce the new slope. The new function is added at transition point but does not contribute to the output prior to that point. The non-linear activation function is responsible for disabling the input prior to and sometimes after the transition points. The following exercise provides a more tangible example.

The exercise will use the network below in a regression problem where you must model a piecewise linear target :



The network has 3 units in the first layer. Each is required to form the target. Unit 0 is pre-programmed and fixed to map the first segment. You will modify weights and biases in unit 1 and 2 to model the 2nd and 3rd segment. The output unit is also fixed and simply sums the outputs of the first layer.

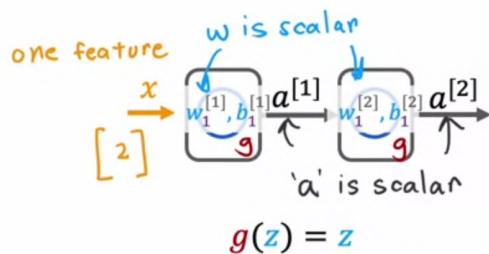
o



▼ Why do activation functions are essential for performance ?

- here $g(z) = z$, which means linear example, after passing a neuron a same impact has been created. So, there is no need of layer 1, which is meaning less. No major learning has been done. **So never use a linear functions within the hidden layer**

Linear Example

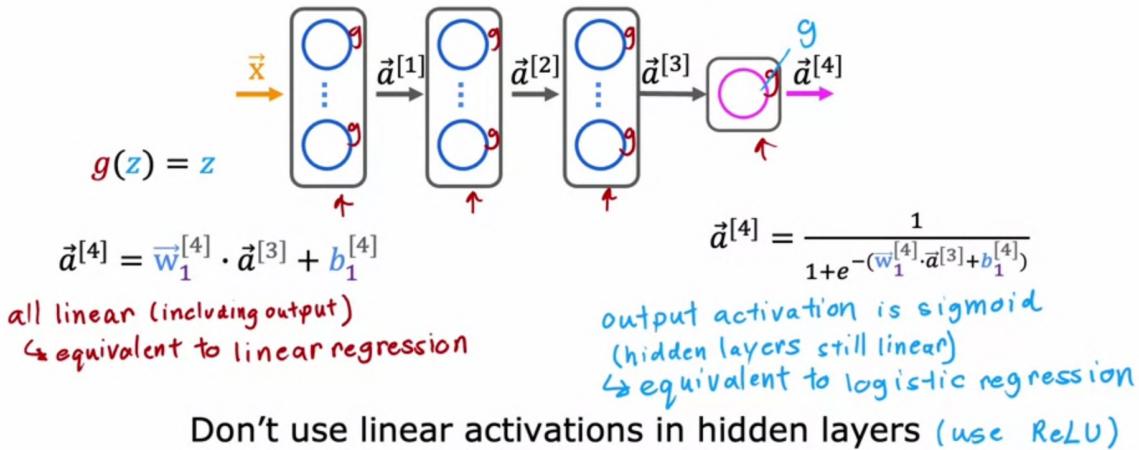


$$\begin{aligned}
 a^{[1]} &= \underbrace{w_1^{[1]} x}_{\downarrow} + b_1^{[1]} \\
 a^{[2]} &= w_1^{[2]} a^{[1]} + b_1^{[2]} \\
 &= w_1^{[2]} (w_1^{[1]} x + b_1^{[1]}) + b_1^{[2]} \\
 \vec{a}^{[2]} &= (\underbrace{\vec{w}_1^{[2]} \vec{w}_1^{[1]}}_{\omega}) x + \underbrace{w_1^{[2]} b_1^{[1]} + b_1^{[2]}}_{b}
 \end{aligned}$$

$$\vec{a}^{[2]} = w x + b$$

$$f(x) = w x + b \text{ linear regression}$$

Example



- ▼ how will you build neural network classification problem with multiple classes ?
- ▼ ReLU often introduces non linearity within the activations. how ?

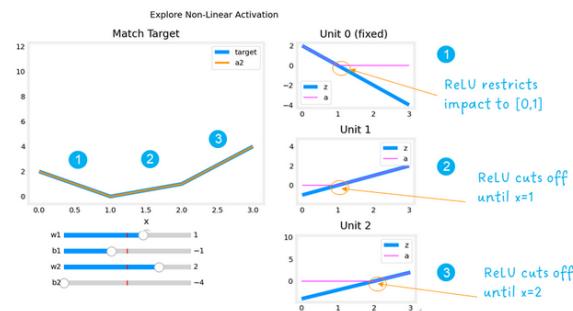
The goal of this exercise is to appreciate how the ReLU's non-linear behavior provides the needed ability to turn functions off until they are needed. Let's see how this worked in this example. The plots on the right contain the output of the units in the first layer.

Starting at the top, unit 0 is responsible for the first segment marked with a 1. Both the linear function z and the function following the ReLU a are shown. You can see that the ReLU cuts off the function after the interval $[0,1]$. This is important as it prevents Unit 0 from interfering with the following segment.

Unit 1 is responsible for the 2nd segment. Here the ReLU kept this unit quiet until after x is 1. Since the first unit is not contributing, the slope for unit 1, $w_1^{[1]}$, is just the slope of the target line. The bias must be adjusted to keep the output negative until x has reached 1. Note how the contribution of Unit 1 extends to the 3rd segment as well.

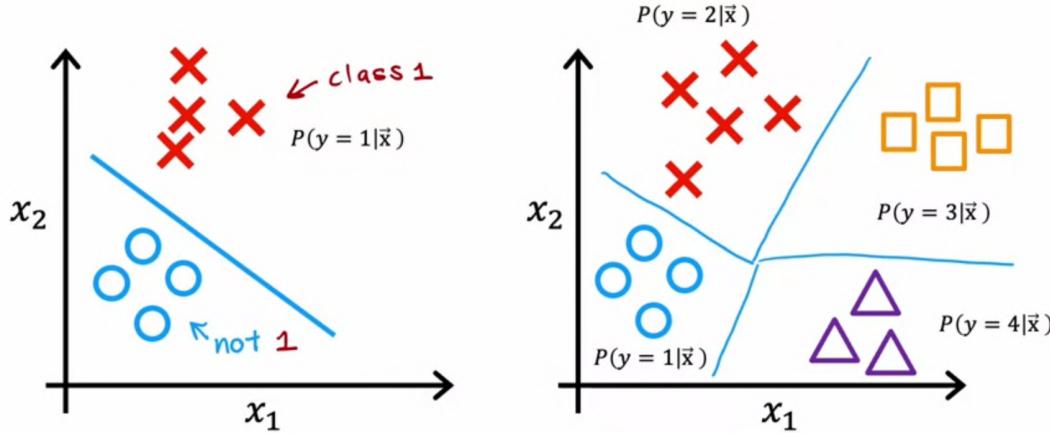
Unit 2 is responsible for the 3rd segment. The ReLU again zeros the output until x reaches the right value. The slope of the unit, $w_2^{[1]}$, must be set so that the sum of unit 1 and 2 have the desired slope. The bias is again adjusted to keep the output negative until x has reached 2.

The "off" or disable feature of the ReLU activation enables models to stitch together linear segments to model complex non-linear functions



- Output labels are more than 2.
- Decision boundaries helps in drawing the classes out of it

Multiclass classification example



- we will be using the SoftMax regression function which is a generalized form of sigmoid function
-

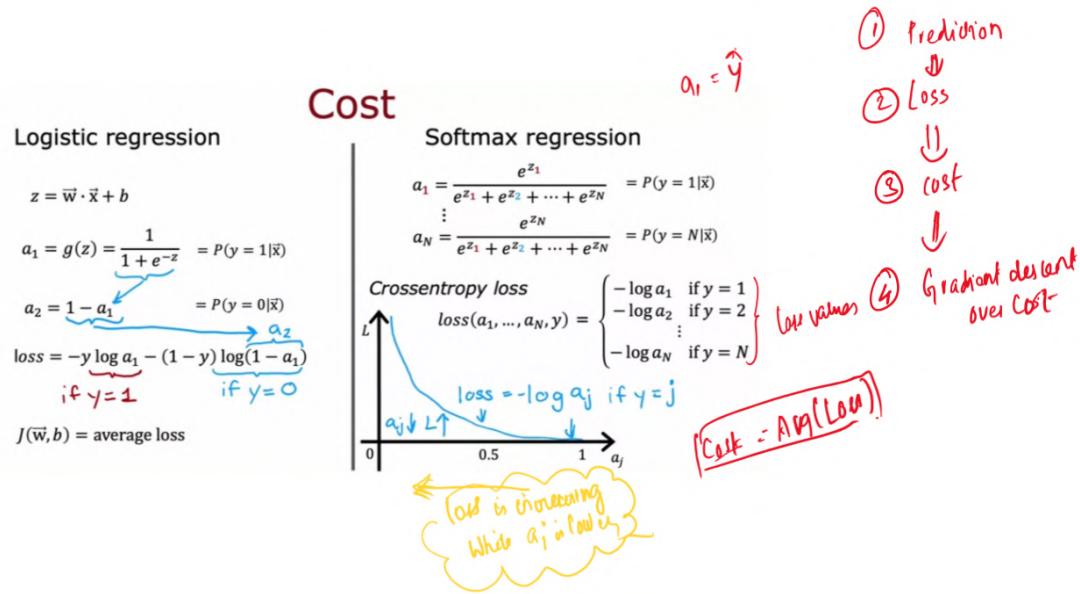
▼ What is SoftMax function ? how it is used in multi-class classification problems?

 - When SoftMax function applied for 2 outputs then it becomes a logistic regression.

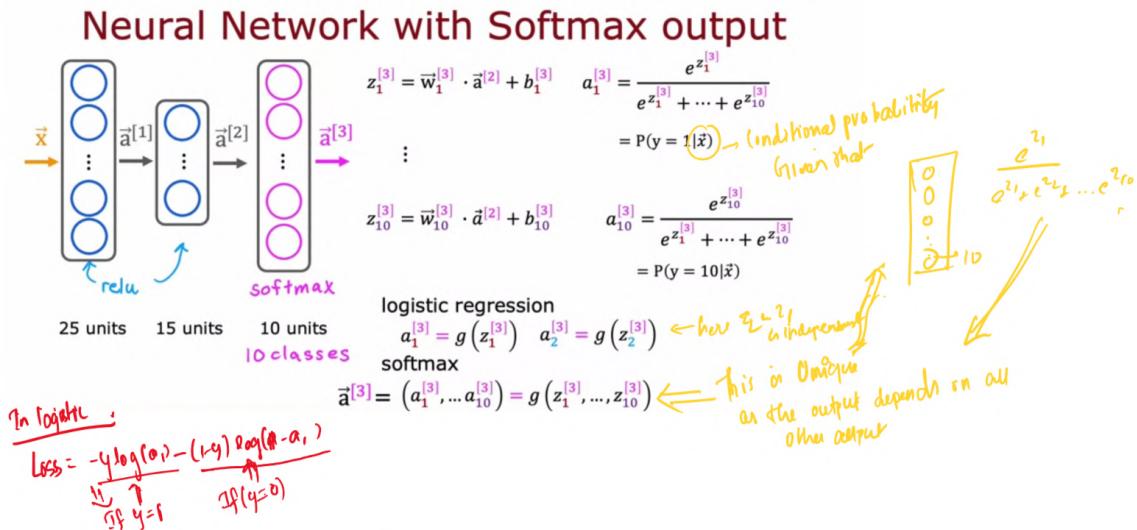
This is probability of being 1

<p>Logistic regression (2 possible output values)</p> $z = \vec{w} \cdot \vec{x} + b$ $\times a_1 = g(z) = \frac{1}{1+e^{-z}} = P(y=1 \vec{x}) \quad 0.11$ $\circ a_2 = 1 - a_1 = P(y=0 \vec{x}) \quad 0.29$ <hr/> <p>Softmax regression (N possible outputs) $y=1, 2, 3, \dots, N$</p> $z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j = 1, \dots, N$ <p>parameters w_1, w_2, \dots, w_N, b_1, b_2, \dots, b_N</p> $a_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j \vec{x})$ <p>note: $a_1 + a_2 + \dots + a_N = 1$</p>	<p>Softmax regression (4 possible outputs) $y=1, 2, 3, 4$</p> $\times z_1 = \vec{w}_1 \cdot \vec{x} + b_1 \quad a_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=1 \vec{x}) \quad 0.30$ $\circ z_2 = \vec{w}_2 \cdot \vec{x} + b_2 \quad a_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=2 \vec{x}) \quad 0.20$ $\square z_3 = \vec{w}_3 \cdot \vec{x} + b_3 \quad a_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=3 \vec{x}) \quad 0.15$ $\triangle z_4 = \vec{w}_4 \cdot \vec{x} + b_4 \quad a_4 = \frac{e^{z_4}}{e^{z_1} + e^{z_2} + e^{z_3} + e^{z_4}} = P(y=4 \vec{x}) \quad 0.35$ <p>for 4 variables</p> <p>output layer</p> <p>Final prediction each class</p>
---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

- cost function within SoftMax function.



- uniqueness of SoftMax function :



- Loss function and steps of implementation :

MNIST with softmax

① specify the model

$$f_{\mathbf{w}, b}(\mathbf{x}) = ?$$

```
import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
from tensorflow.keras.losses import
SparseCategoricalCrossentropy()
model.compile(loss= SparseCategoricalCrossentropy() )
model.fit(X,Y,epochs=100)
Note: better (recommended) version later.
Don't use the version shown here!
```

② specify loss and cost

$$L(f_{\mathbf{w}, b}(\mathbf{x}), \mathbf{y})$$

③ Train on data to minimize $J(\mathbf{w}, b)$

Number $\{1, b\}$ predicted
 $[1, 2, 3, 4, 5, 6, 7, 8, 9]$
 \downarrow
output looks like $(\hat{y} = 2)$
 $[0, 1, 0, 0, 1, 0, 0, 0, 0, 0]$
sparse categorical (a_i)
categorical entropy

Binary classes
 $[1, 0]$
 \downarrow
Binary cross entropy

▼ how to implement the improved version of SoftMax function ?

- There occurs a rounding off error in the above implementation of SoftMax function
- The memory of each digits are finite, so instead of storing them in separate variable, it would be optimal for us to express within the equation for better and accurate(re-arranges to a optimized execution) results.

•

Numerical Roundoff Errors

More numerically accurate implementation of logistic loss:

Logistic regression:

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

Original loss

$$loss = -y \log(\hat{a}) - (1-y) \log(1-\hat{a})$$

More accurate loss (in code)

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

$$1 + \frac{1}{10,000} \quad 1 - \frac{1}{10,000}$$

model = Sequential([
 Dense(units=25, activation='relu'),
 Dense(units=15, activation='relu'),
 Dense(units=1, activation='sigmoid')
])

model.compile(loss=BinaryCrossEntropy())

model.compile(loss=BinaryCrossEntropy(from_logits=True))

Might Not be Major in Sigmoid Function
has Major Impact on "Softmax"

Instead of Using
Sigmoid
Linear + (Softmax)

To avoid numerical
errors caused due
to Memory allocation.

More numerically accurate implementation of softmax

Softmax regression

$$(a_1, \dots, a_{10}) = g(z_1, \dots, z_{10})$$

$$\text{Loss} = L(\vec{a}, y) = \begin{cases} -\log a_1 & \text{if } y = 1 \\ -\log a_{10} & \text{if } y = 10 \end{cases}$$

```
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='softmax')
])
model.compile(loss=SparseCategoricalCrossentropy())
```

'linear'

More Accurate

$$L(\vec{a}, y) = \begin{cases} -\log \frac{e^{z_1}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 1 \\ -\log \frac{e^{z_{10}}}{e^{z_1} + \dots + e^{z_{10}}} & \text{if } y = 10 \end{cases}$$

```
model.compile(loss=SparseCategoricalCrossEntropy(from_logits=True))
```



- Implementation :

- More accurate implementation of Softmax function

MNIST (more numerically accurate)

```
model import tensorflow as tf
from tensorflow.keras import Sequential
from tensorflow.keras.layers import Dense
model = Sequential([
    Dense(units=25, activation='relu'),
    Dense(units=15, activation='relu'),
    Dense(units=10, activation='linear') ])
loss from tensorflow.keras.losses import
      SparseCategoricalCrossentropy
model.compile(..., loss=SparseCategoricalCrossentropy(from_logits=True))
fit model.fit(X, Y, epochs=100)
predict logits = model(X) ← not a1...a10
                     is z1...z10
f_x = tf.nn.softmax(logits)
```

- More accurate way of implementing logistic function :

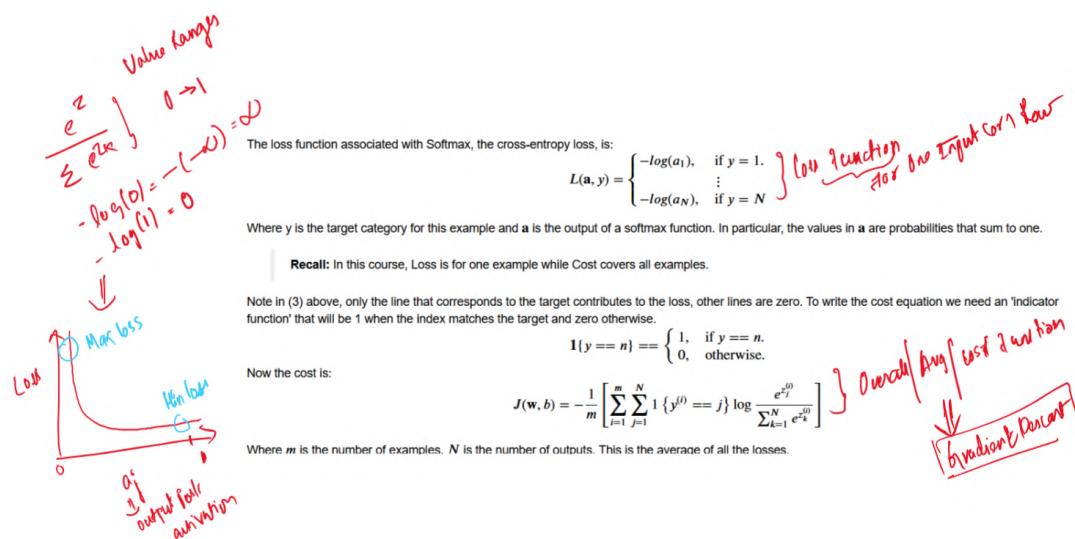
logistic regression (more numerically accurate)

```

model = Sequential([
    Dense(units=25, activation='sigmoid'),
    Dense(units=15, activation='sigmoid'),
    Dense(units=1, activation='linear')
])
from tensorflow.keras.losses import
    BinaryCrossentropy
model.compile(..., BinaryCrossentropy(from_logits=True))
model.fit(X,Y,epochs=100)
fit      logit = model(X)    ↗
predict  f_x = tf.nn.sigmoid(logit)  ↗

```

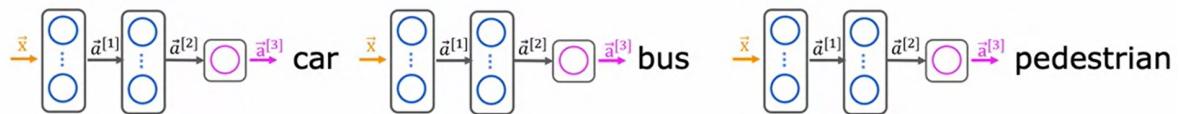
▼ Loss and cost function of SoftMax regression function ?



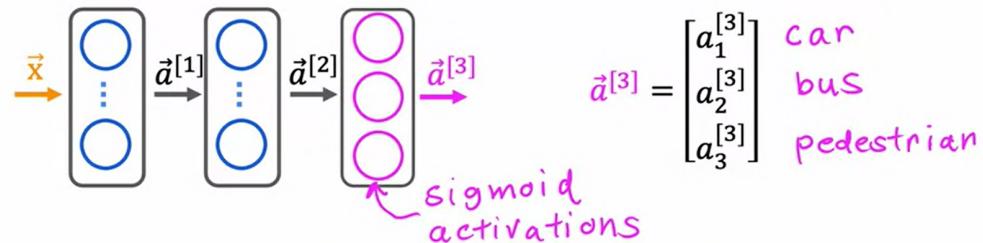
▼ What is multi - label classification ?

- In image recognition, sometime we predict more than one object: eg. Detecting traffic signal and pedestrians simultaneously. These kind of problems where the detection or outcome required is more than one : is called multi-label classification

Multi-label Classification



Alternatively, train one neural network with three outputs



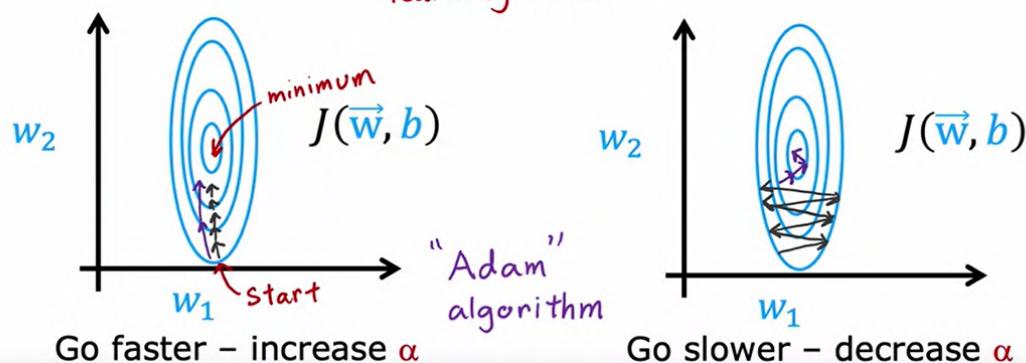
▼ What are the other methods than reducing cost function by gradient descent ?

- When applying gradient descent the weights tend to converge when the cost function is minimum, if the learning rate is too low then this would take smaller steps though its converging → so to make this faster we are using an optimizer called "**Adam**"

Gradient Descent

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

learning rate



- This Adam optimizer guides and changes the learning rate such that to converge faster by manipulating the alpha.
- The learning is not global here.

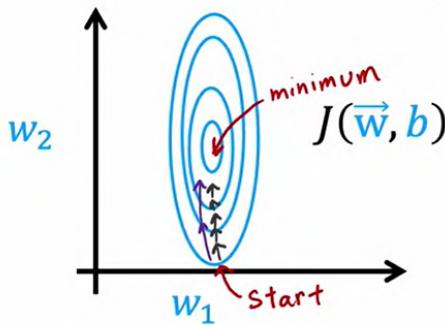
Adam Algorithm Intuition

Adam: Adaptive Moment estimation *not just one α*

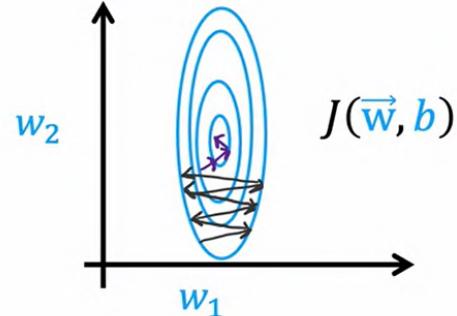
$$\begin{aligned} w_1 &= w_1 - \underbrace{\alpha_1}_{\vdots} \frac{\partial}{\partial w_1} J(\vec{w}, b) \\ w_{10} &= w_{10} - \underbrace{\alpha_{10}}_{\vdots} \frac{\partial}{\partial w_{10}} J(\vec{w}, b) \\ b &= b - \underbrace{\alpha_{11}}_{\vdots} \frac{\partial}{\partial b} J(\vec{w}, b) \end{aligned}$$

- here the alpha changes with the nature of change of parameters w and b.
Below picture depicts it :

Adam Algorithm Intuition



If w_j (or b) keeps moving
in same direction,
increase α_j .



If w_j (or b) keeps oscillating,
reduce α_j .

- its robust with the learning rate optimization, the implementation code can be found below :

MNIST Adam model

```
model = Sequential([
    tf.keras.layers.Dense(units=25, activation='sigmoid'),
    tf.keras.layers.Dense(units=15, activation='sigmoid'),
    tf.keras.layers.Dense(units=10, activation='linear')
])
```

compile

$$\alpha = 10^{-3} = 0.001$$

```
model.compile(optimizer=tf.keras.optimizers.Adam(learning_rate=1e-3),
              loss=tf.keras.losses.SparseCategoricalCrossentropy(from_logits=True))
```

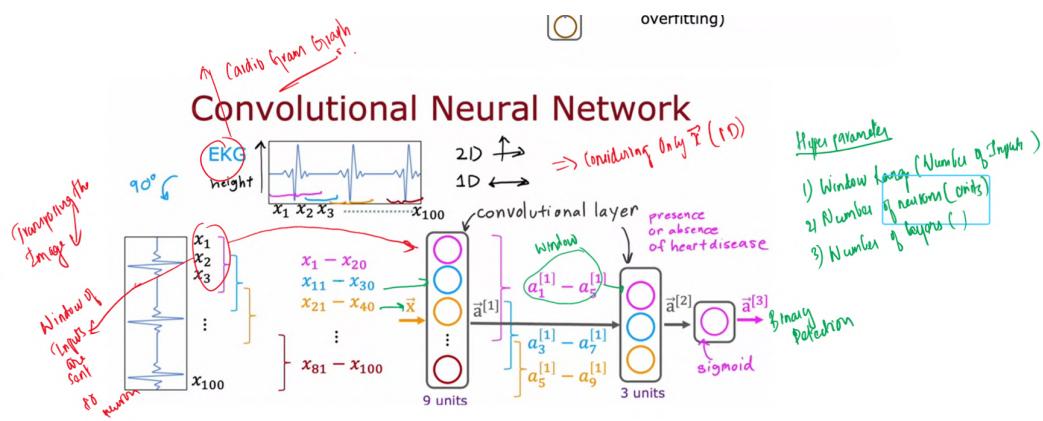
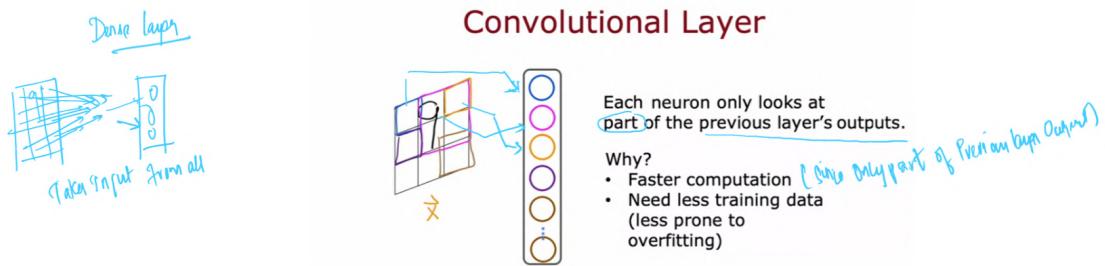
fit

```
model.fit(X, Y, epochs=100)
```

▼ What are the alternative way or types ?

▼ Different type of neural network

- Convolutional layer :

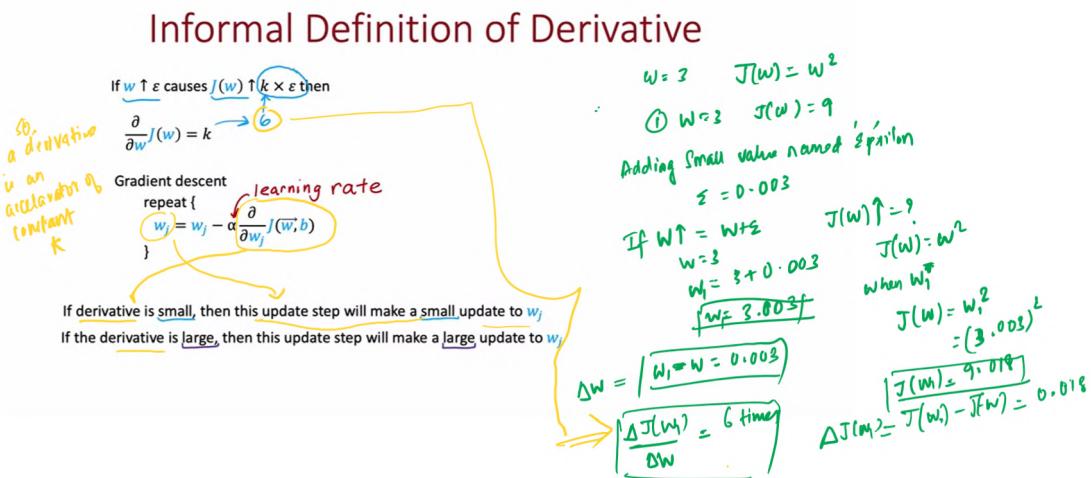


- you can see a sliding window which determines the number of inputs to be sent to other layer, instead of considering all the inputs.

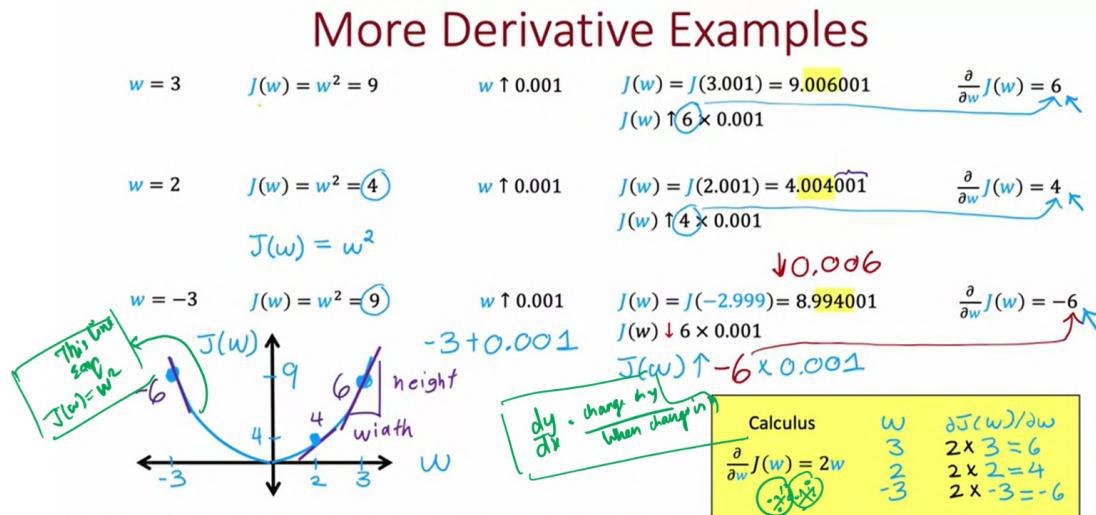
-

▼ How the back propagation works ?

▼ Understanding derivative:



- The derivative is impacted by the weights of the parameter not Loss function :



- sympy : library is used to implement the derivative of specific function.

- o

Even More Derivative Examples

$$\begin{array}{llll}
 w = 2 & \left\{ \begin{array}{lll}
 J(w) = w^2 = 4 & \frac{\partial}{\partial w} J(w) = 2w = 4 & w \uparrow \frac{0.001}{\epsilon} \\
 J(w) = w^3 = 8 & \frac{\partial}{\partial w} J(w) = 3w^2 = 12 & w \uparrow \epsilon \\
 J(w) = w = 2 & \frac{\partial}{\partial w} J(w) = 1 & w \uparrow \epsilon \\
 J(w) = \frac{1}{w} = 0.5 & \frac{\partial}{\partial w} J(w) = -\frac{1}{w^2} = -\frac{1}{4} & w \uparrow \epsilon \\
 \end{array} \right. & \begin{array}{ll}
 J(w) = 4.004001 & J(w) \uparrow 4 \times \epsilon \\
 J(w) = 8.012006 & J(w) \uparrow 12 \times \epsilon \\
 J(w) = 2.001 & J(w) \uparrow 1 \times \epsilon \\
 \cancel{J(w) = 0.5 - 0.00025} & \cancel{J(w) = 0.49975} \\
 J(w) \uparrow \frac{1}{2.001} & J(w) \uparrow -\frac{1}{4} \times \epsilon
 \end{array} \\
 \boxed{\frac{\partial}{\partial w} J(w) \quad w \uparrow \epsilon \quad J(w) \uparrow k \times \epsilon} &
 \end{array}$$

A note on derivative notation

If $J(w)$ is a function of one variable (w),

$$d \quad \frac{d}{dw} J(w)$$

If $J(w_1, w_2, \dots, w_n)$ is a function of more than one variable,

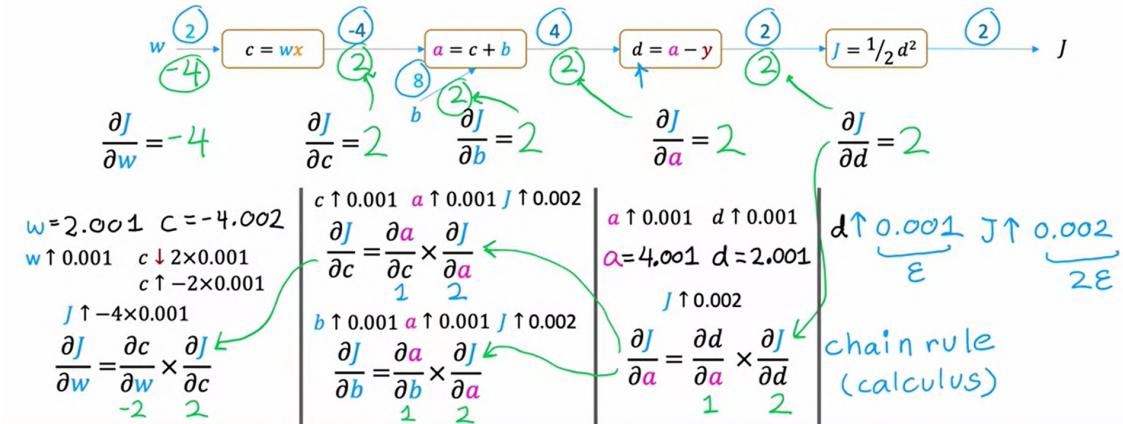
$$\partial \quad \frac{\partial}{\partial w_i} J(w_1, w_2, \dots, w_n)$$

▼ How backpropagation works ?

- the application of partial derivative from left to right helps in normalizing the weights of the neuron. The partial derivative applies chain rule as well.

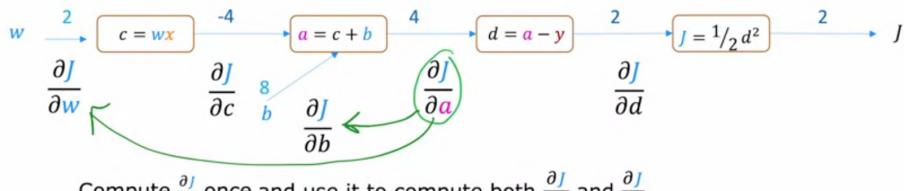
Computing the Derivatives

$$w = 2 \quad b = 8 \quad x = -2 \quad y = 2 \quad a = wx + b \quad J = \frac{1}{2}(a - y)^2$$



- why back propagation is essential ?
 - Right to left propagation essentially helps to understand, how the output is impacted by the input. How they are interlinked.
 - This back propagation helps in achieving highly efficient computing by minimizing the parameters mentioned in the image below :

Backprop is an efficient way to compute derivatives



Compute $\frac{\partial J}{\partial a}$ once and use it to compute both $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial b}$.

If N nodes and P parameters, compute derivatives

in roughly $N + P$ steps rather than $N \times P$ steps.

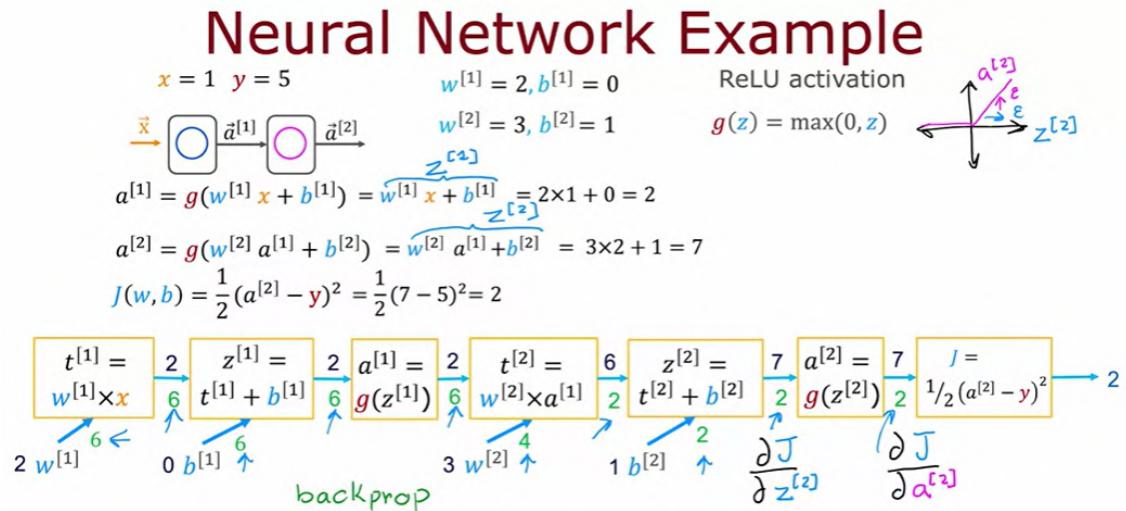
N	P	N+P	NXP
10,000	100,000	1.1×10^5	10^9

- N : number of nodes → P : Number of parameters | $N+P =$ total number of parameters used.

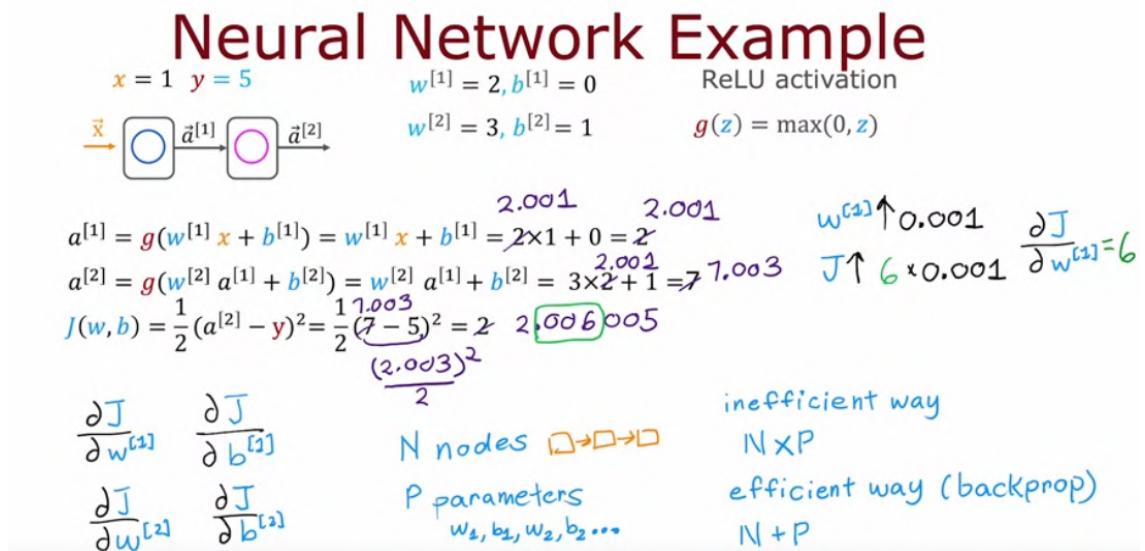
- Left → right : forward propagation is used for calculating the cost function and back propagation are used to calculate the derivatives (parameters of the weight.).

▼ How back propagation works in larger neural network ?

- implementation diagram :



- for N nodes and p parameters if only forward propagations exists then the adjustments would take N times P, however the backpropagation helps it to reduce to N + P adjustments. Earlier times these derivatives are calculated manually.



▼ How to apply machine learning ? with diagnosis

▼ What if you incur with higher losses and undesirable results ?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- adding polynomial features
- Increasing and decreasing Lambda

▼ How to evaluate and choose the model ?

- when we have one feature (meaning x_1 and y) then we can plot and diagnose the model, if it is more than 2 then we invariably need an evaluation metric to diagnose the model. because now it has become a multidimensional model
- Splitting the data into train and test data set, fit the model on training data and analyze the prediction data on test data. Sometimes it makes sense to test the prediction on training dataset.

Train/test procedure for linear regression (with squared error cost)

Fit parameters by minimizing cost function $J(\vec{w}, b)$

$$\rightarrow J(\vec{w}, b) = \left[\frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m_{train}} \sum_{j=1}^n w_j^2 \right]$$

Compute test error:

$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}} \left[\sum_{i=1}^{m_{test}} (f_{\vec{w}, b}(\vec{x}_{test}^{(i)}) - y_{test}^{(i)})^2 \right] \quad \text{with } \sum_{j=1}^n w_j^2$$

Compute training error:

$$J_{train}(\vec{w}, b) = \frac{1}{2m_{train}} \left[\sum_{i=1}^{m_{train}} (f_{\vec{w}, b}(\vec{x}_{train}^{(i)}) - y_{train}^{(i)})^2 \right]$$

Train/test procedure for classification problem

0 / 1

Fit parameters by minimizing $J(\vec{w}, b)$ to find \vec{w}, b

E.g.,

$$J(\vec{w}, b) = -\frac{1}{m_{train}} \sum_{i=1}^{m_{train}} \underbrace{\left[y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right]}_{+ \frac{\lambda}{2m_{train}} \sum_{j=1}^n w_j^2}$$

Compute test error:

$$J_{test}(\vec{w}, b) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \underbrace{\left[y_{test}^{(i)} \log(f_{\vec{w}, b}(\vec{x}_{test}^{(i)})) + (1 - y_{test}^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}_{test}^{(i)})) \right]}$$

Compute train error:

$$J_{train}(\vec{w}, b) = -\frac{1}{m_{train}} \sum_{i=1}^{m_{train}} \left[y_{train}^{(i)} \log(f_{\vec{w}, b}(\vec{x}_{train}^{(i)})) + (1 - y_{train}^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}_{train}^{(i)})) \right]$$

Train/test procedure for classification problem

fraction of the test set and the fraction of the train set that the algorithm has misclassified.

count $\hat{y} \neq y$

$$\hat{y} = \begin{cases} 1 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) \geq 0.5 \\ 0 & \text{if } f_{\vec{w}, b}(\vec{x}^{(i)}) < 0.5 \end{cases}$$

$J_{test}(\vec{w}, b)$ is the fraction of the test set that has been misclassified.

$J_{train}(\vec{w}, b)$ is the fraction of the train set that has been misclassified.

▼ How the training/test/validation dataset works ? why is it needed ?

- Features of good model
 - the test misclassification is lower and loss function is minimum
 - Difference between the train and test error will be lower.
 - They perform good on new dataset and
- There is flaw when you identify simply a model with minimum job cost function when you are iterating through multiple polynomials :

Model selection (choosing a model)

- $d=1$ 1. $f_{\bar{w},b}(\vec{x}) = w_1x + b \rightarrow w^{<1>}, b^{<1>} \rightarrow J_{test}(w^{<1>}, b^{<1>})$
- $d=2$ 2. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + b \rightarrow w^{<2>}, b^{<2>} \rightarrow J_{test}(w^{<2>}, b^{<2>})$
- $d=3$ 3. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + w_3x^3 + b \rightarrow w^{<3>}, b^{<3>} \rightarrow J_{test}(w^{<3>}, b^{<3>})$
- \vdots
- $d=10$ 10. $f_{\bar{w},b}(\vec{x}) = w_1x + w_2x^2 + \dots + w_{10}x^{10} + b \rightarrow J_{test}(w^{<10>}, b^{<10>})$
- Choose $w_1x + \dots + w_5x^5 + b \quad d=5 \quad J_{test}(w^{<5>}, b^{<5>})$

How well does the model perform? Report test set error $J_{test}(w^{<5>}, b^{<5>})$?

The problem: $J_{test}(w^{<5>}, b^{<5>})$ is likely to be an optimistic estimate of generalization error (ie. $J_{test}(w^{<5>}, b^{<5>}) <$ generalization error). Because an extra parameter d (degree of polynomial) was chosen using the test set.

w, b are overly optimistic estimate of generalization error on training data.

- So we are going to implement the validation dataset as a solution :

		Training/cross validation/test set	
size	price	validation set	development set
2104	400		
1600	330		
2400	369		
1416	232		
3000	540		
1985	300		
1534	315	training set 60%	$(x^{(1)}, y^{(1)})$ \vdots $(x^{(m_{train})}, y^{(m_{train})})$
1427	199		
1380	212	cross validation 20%	$(x_{cv}^{(1)}, y_{cv}^{(1)})$ \vdots $(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
1494	243	test set 20%	$(x_{test}^{(1)}, y_{test}^{(1)})$ \vdots $(x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$M_{train} = 6$

$M_{cv} = 2$

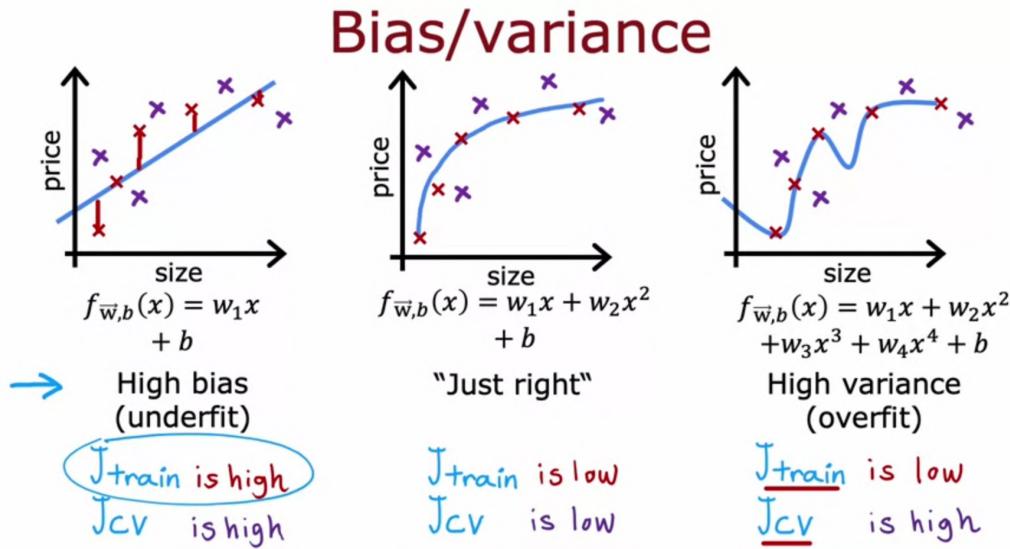
$M_{test} = 2$

- We will be choosing the model which has lower cost function in the validation dataset, then we apply the model to predict on testing data.
- Why do we need validation dataset ?
 - Validation dataset is used to identify the fraction of misclassification that happens within the model.

- Model is not exposed to test dataset → so we can use test dataset to estimate the generalization power of the model.

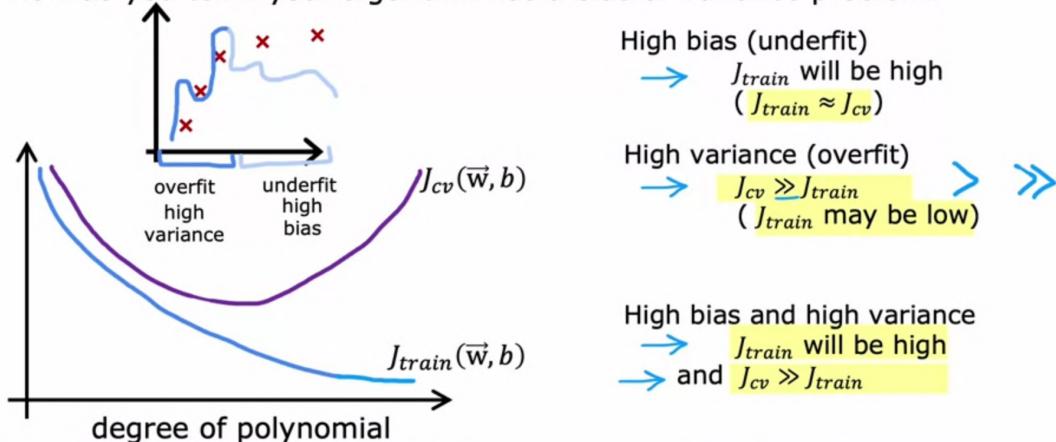
▼ What is bias and variance ? and how to diagnose it

- Bias : does not even predict well on the training dataset, variance : even if the training dataset performs well → new dataset is unpredictable.



Diagnosing bias and variance

How do you tell if your algorithm has a bias or variance problem?



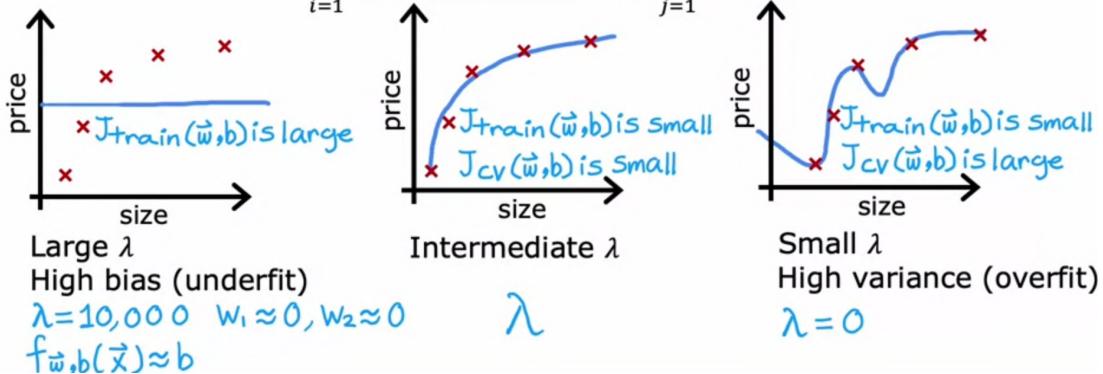
▼ How regularization helps in addressing the problems caused by bias and variance

- how regularization impacts the algorithm ? :

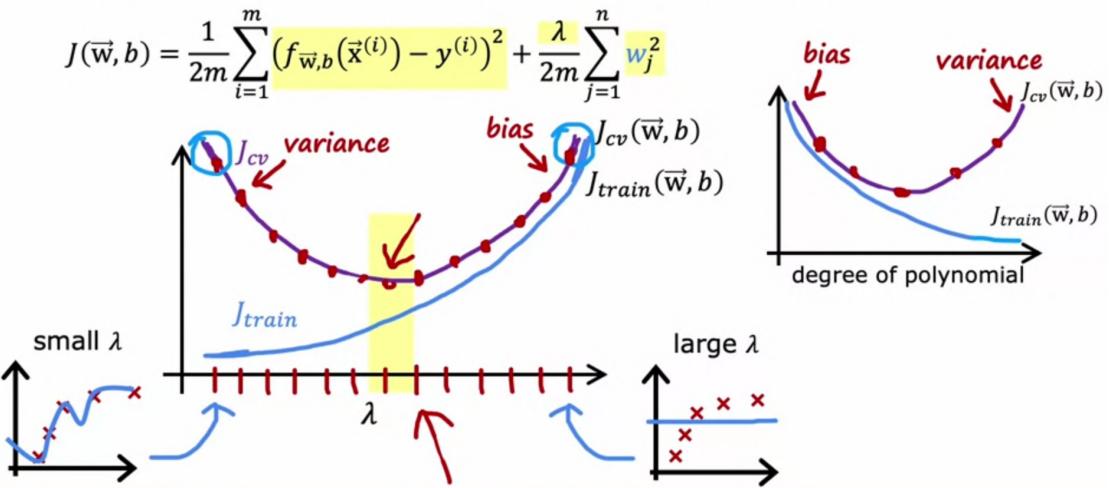
Linear regression with regularization

Model: $f_{\vec{w}, b}(x) = \underline{w_1}x + \underline{w_2}x^2 + \underline{w_3}x^3 + \underline{w_4}x^4 + b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$



Bias and variance as a function of regularization parameter λ



- Usually the training performance are expected to be lower than human level performance and we can use the human level performance as a benchmark for starting.
- Most of the cases, the training error should be less than the validation error. High training error would the bias is high.
- What is the level of error you can reasonably hope to get to ?

- Human level performance : how much error % is to be expected from a human.
- Competing with the algorithms performance
- Guessing based on subject expertise.

▼ How to set benchmark the performance of model ?

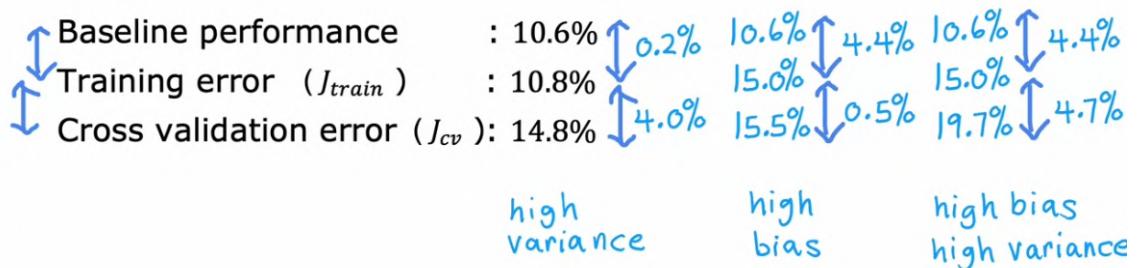
▼ Benchmark the model performance

- Check the human level performance
- prior experience
- maximum performance gotten by competing algorithm on same applications
- Getting to know others performance

▼ How to evaluate the performance

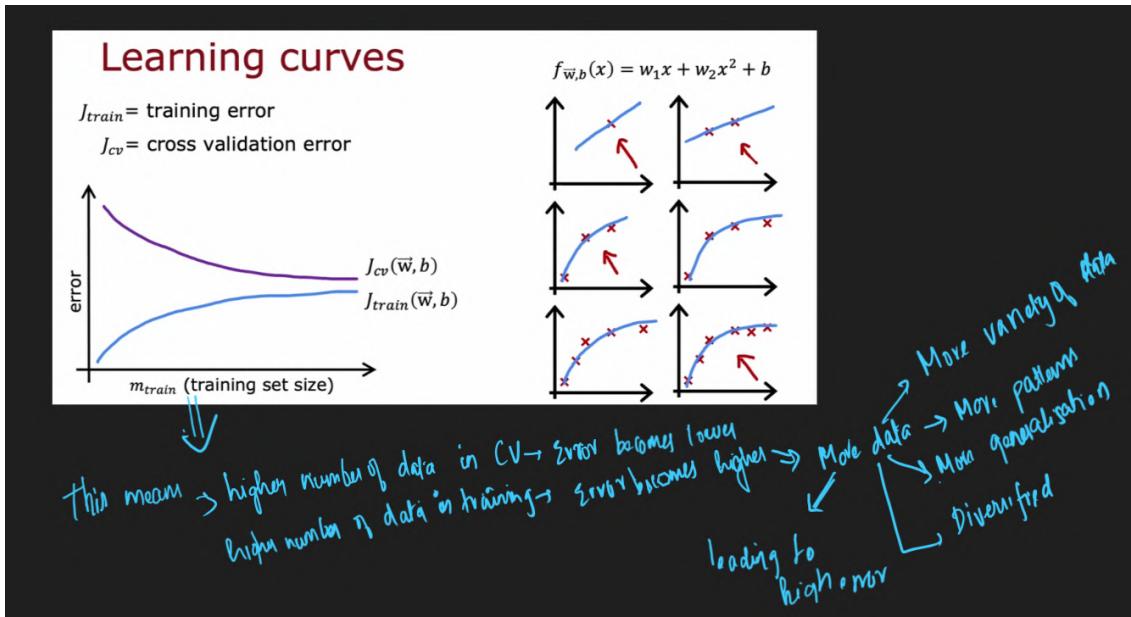
- Difference between the baseline (benchmark : like human level performance or any competing algorithm) and validation error explains the bias
- Difference between the validation and test explains the variance
- Thus we can address the overfitting and underfitting issue.

Bias/variance examples

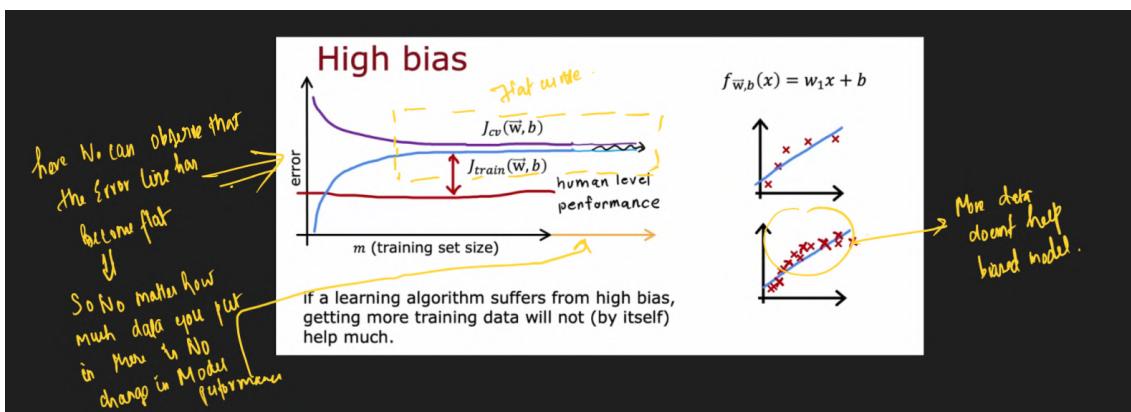


▼ How to learning curve to establish an baseline model ?

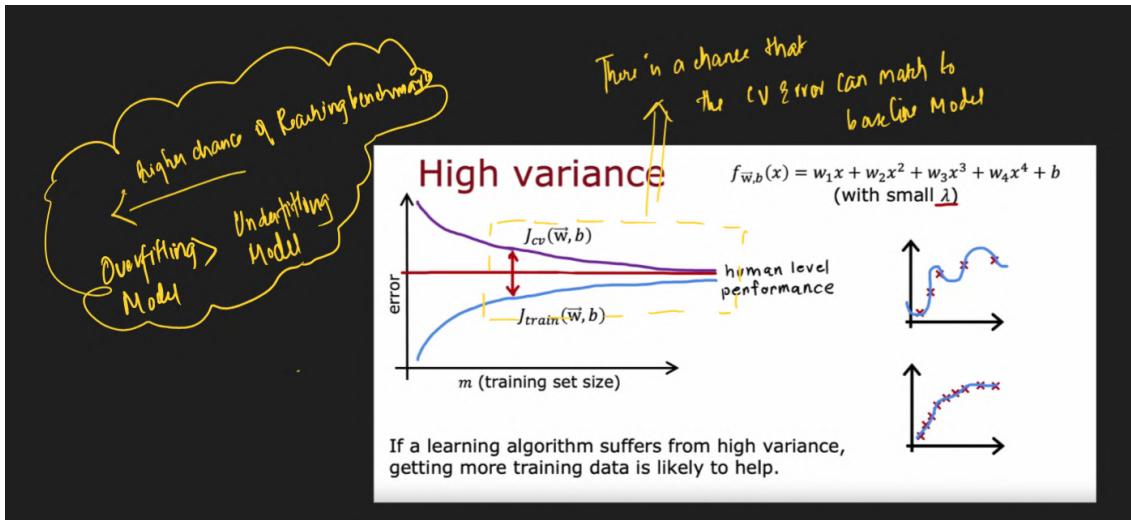
-



- When the model is biased as seen below : no matter how large data imputed into the model it is almost not possible to improve its performance, hence when we identify bias with baseline performance then it should be addressed first. Also ensure to include diversified dataset points into the model.



- High variance model (mentioned below): this model on exposing to more and more data, there is probability where can converge with baseline model. Hence the more prone model is the biased model where adding data is not sufficient to address the issue.



▼ How to proceed further while encountering the bias and variance ?

Overfitting problem	Underfitting problem
Get more training examples	Try adding more features
Try increasing lambda : more regularization	Try adding more polynomial features
Try smaller sets of features	Try decreasing the lambda : fewer regularization is achieved

•

Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

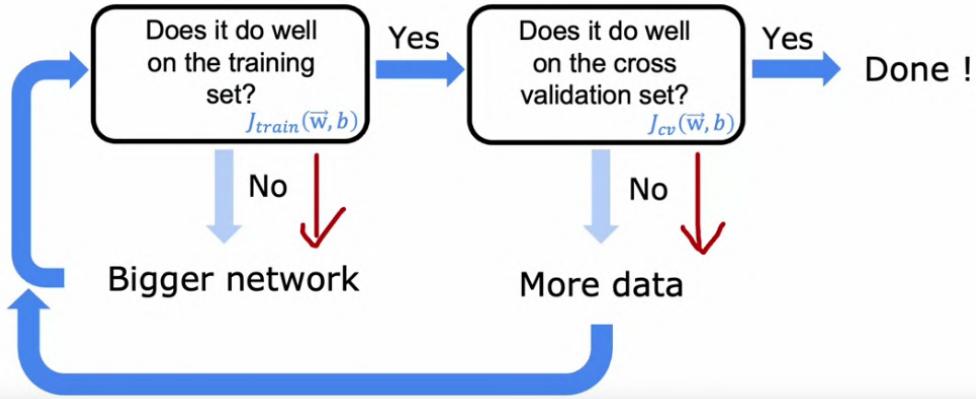
- Get more training examples fixes high variance
- Try smaller sets of features ~~x, x², ..., xⁿ~~ fixes high variance
- Try getting additional features ← fixes high bias
- Try adding polynomial features (x₁², x₂², x₁x₂, etc) ← fixes high bias
- Try decreasing λ ← fixes high bias
- Try increasing λ ← fixes high variance

▼ How the bias and variance concept applied to neural network ?

- Large neural network has lower bias as they may be prone to overfitting.
- Very large neural network could be computationally expensive.
 - Cons
 - More data can be a limiting factor
 - Needs GPU and other accelerators to run the larger neural network

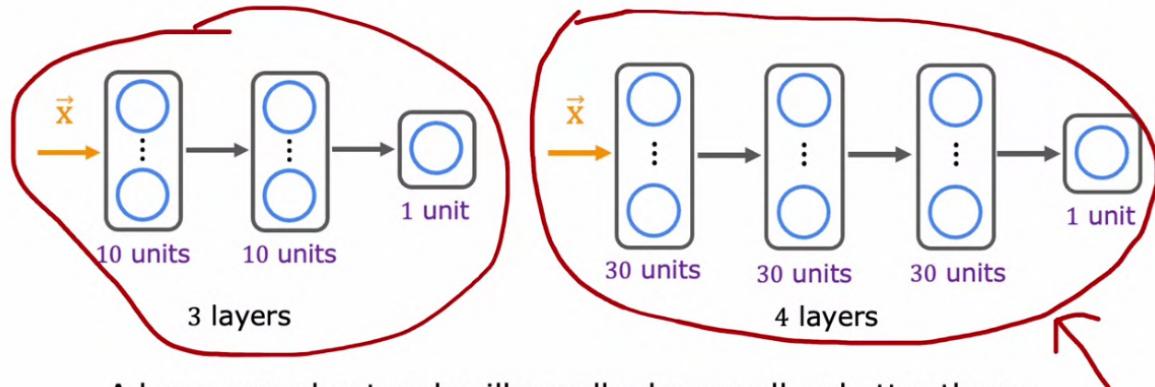
Neural networks and bias variance

Large neural networks are low bias machines



- A idle process would be get a overfitting model with larger neural network and exposing them to new data as well as regularize them to achieve human level data.

Neural networks and regularization



A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

- How regularization applied to neural network applied below ?

Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{x}^{(i)}), y^{(i)}) + \frac{\lambda}{2m} \sum_{\text{all weights } \mathbf{W}} (\mathbf{w}^2)$$

Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

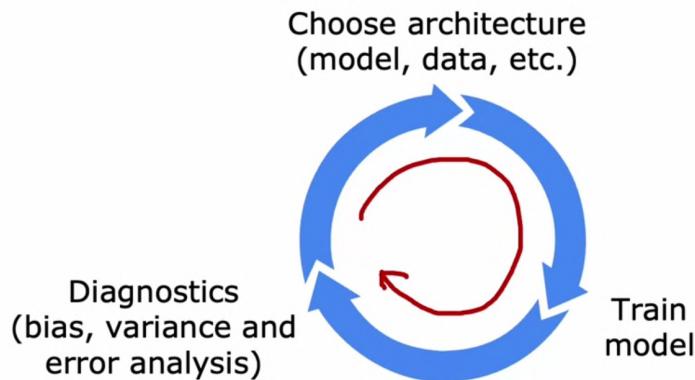
Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

- ▼ How to develop a machine learning system ?

- Iterative loop of development.

Iterative loop of ML development



▼ Error Analysis : How does it look like ?

- Manual diagnosing the errors will be small and tiny impactful activity compared to automatic and computers way.
- Applicable where the bench mark is human level performance meaning : applicable to tasks that human can do better. Prediction may seem better with machine rather than human

▼ How to overcome the problems in machine learning :

- Adding data : on error analysis, if any specific kind of information is misclassified then the specific type of data can be added for betterment of the model.
- Data augmentation for image classification and audio recognition for adding data : new training examples are created by performing few minor modification(rotation , mirroring, introducing new lines within the image) to generate new training data set. Adding more background noise to the existing voice : Data augmentation.
- Never add random noise or meaningless noise.

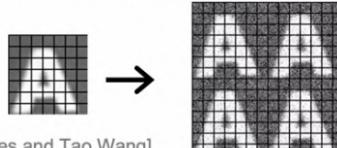
Data augmentation by introducing distortions

Distortion introduced should be representation of the type of noise/distortions in the test set.



Audio:
Background noise,
bad cellphone connection

Usually does not help to add purely random/meaningless noise to your data.



x_i =intensity (brightness) of pixel i
 $x_i \leftarrow x_i + \text{random noise}$

[Adam Coates and Tao Wang]

- Data Synthesis(higher application in computer vision) : Creating brand new examples or training data meaning existing training data is untouched. Manually can be created by work or Image generating LLM's could be employed here.
- **Our approach must foster : Data centric approach as the existing codes are well written and needs higher expertise in improving its approach**

Engineering the data used by your system

Conventional
model-centric
approach:

$$\text{AI} = \text{Code} + \text{Data}$$

(algorithm/model)

work on this

Data-centric
approach:

$$\text{AI} = \text{Code} + \text{Data}$$

(algorithm/model)

work on this

- ▼ What is transfer learning , why should be augmented into our existing approach ?

- Cannot be applied to all, however if possible then it could be powerful

Transfer learning summary

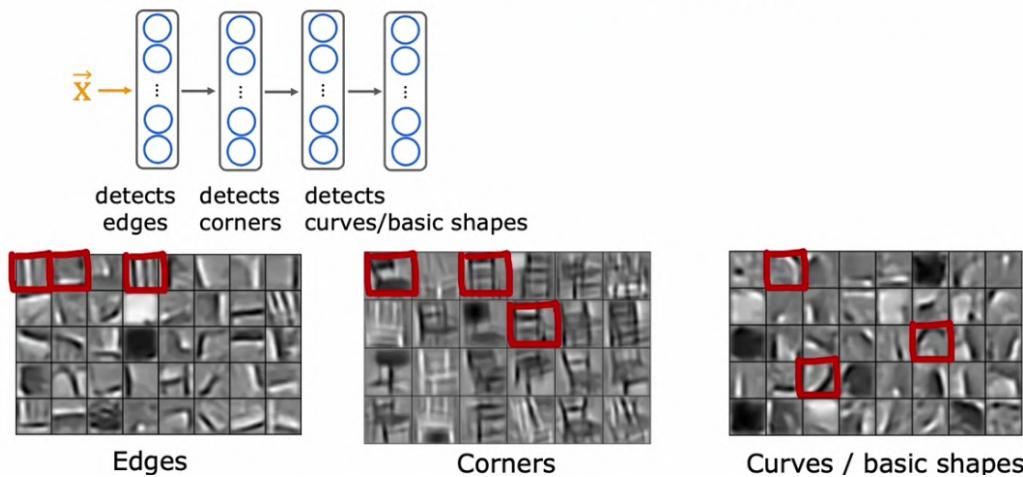
1. Download neural network parameters pretrained on a large dataset with same input type (e.g., images, audio, text) as your application (or train your own). *1 million images*
2. Further train (fine tune) the network on your own data.

1000 images

50 images

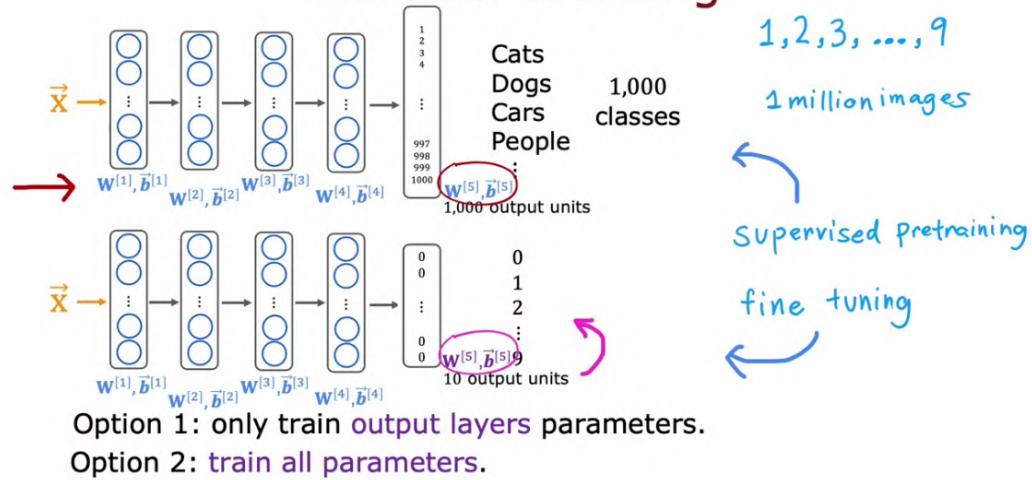
- When we don't have much of training data or when we are in need of adding more data.
- Transfer learning allows to download and use the best model available on internet as they might have already researched and proved to provide better results for any application.
- The neurons with neural layer tries to learn the detection skills rather than detecting output hence, this allows the user to use the same model for different use cases of the same application(image NN on Image data not on audio data)

Why does transfer learning work?



- In transfer learning the output layer is replaced rather than adding a layer
 - reason for not adding a layer
 - Computationally heavy
 - Net impact would be lower
- Options in transfer learning
 - option 1 (Supervised pre training): training only the output layer, meaning the weights and biases of previous layers remains untouched.
 - option 2 (fine-tuning): training all parameters > training all weights and biases of all layer.

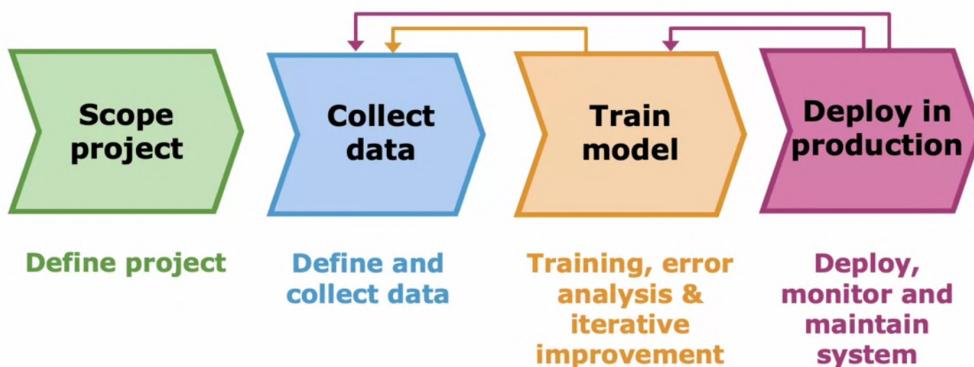
Transfer learning



- ▼ How full cycle of a machine learning project would look like ?

-

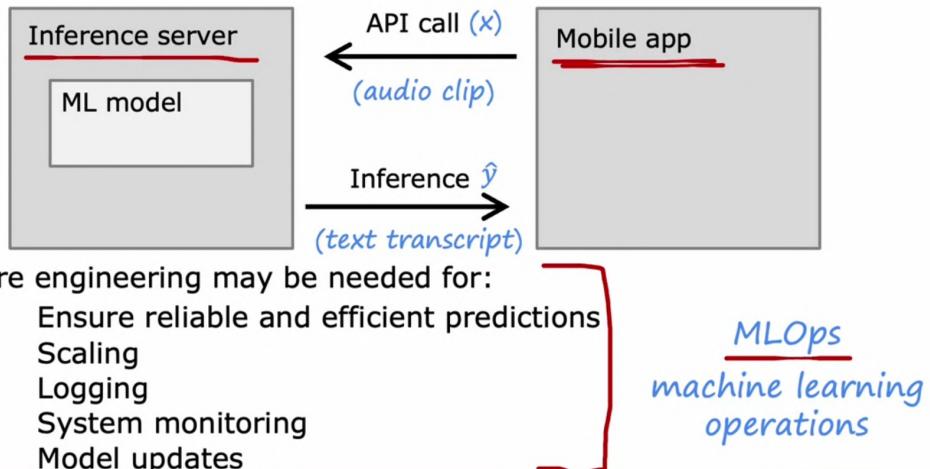
Full cycle of a machine learning project



→ Few software requirements for deployment mentioned below.

-

Deployment



▼ Ethics in ML

-

Bias

Hiring tool that discriminates against women.

Facial recognition system matching dark skinned individuals to criminal mugshots.

Biased bank loan approvals.

Toxic effect of reinforcing negative stereotypes.

Adverse use cases

Deepfakes

Spreading toxic/incendiary speech through optimizing for engagement.

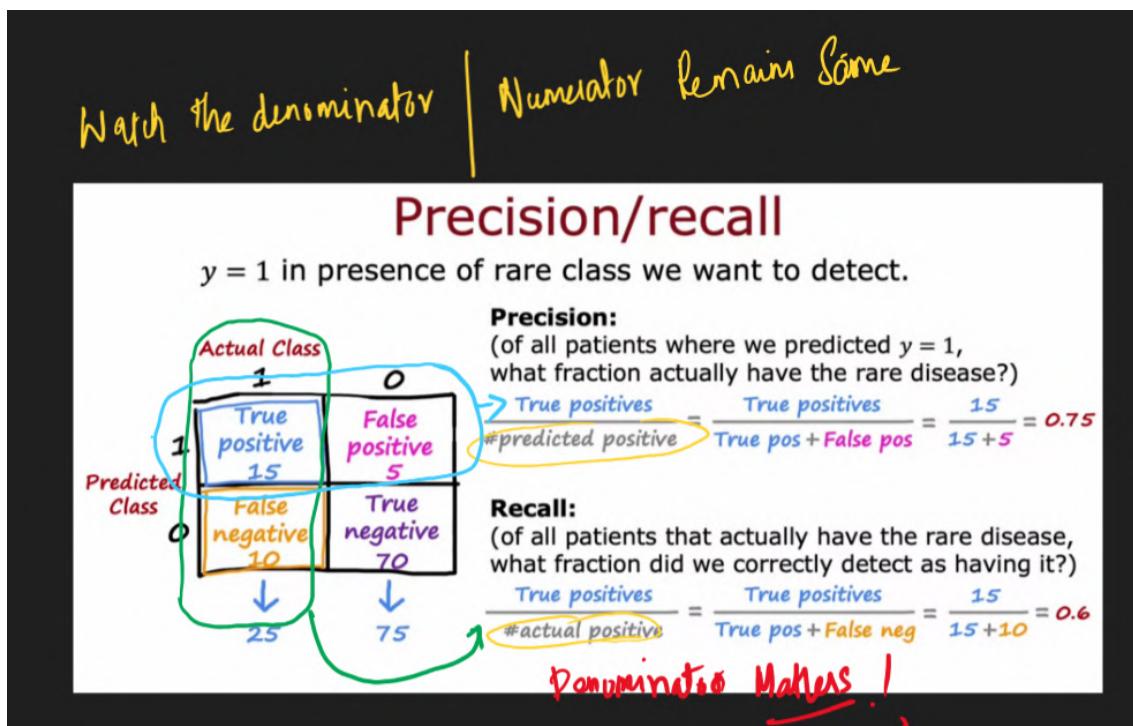
Generating fake content for commercial or political purposes.

Using ML to build harmful products, commit fraud etc.

Spam vs anti-spam : fraud vs anti-fraud.

▼ How to measure Error Metrics for skewed datasets ?

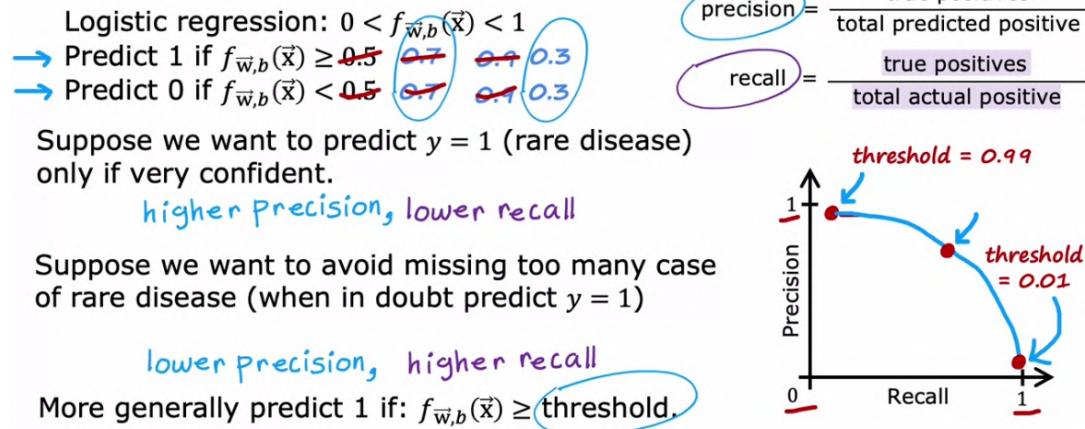
- Precision/Recall : In precision(speaks about model), if a person is diagnosed as sick there 75% chance that he is sick. In recall(speaks about population), the model is capable to predict 60% of total sick persons.



- High precision & High recall is not always possible.

- If the threshold for classification is increased which means that precision becomes higher, however this lowers the recall which means most of the patients are left out instead of identification.

Trading off precision and recall



- It is not always possible to manipulate the threshold as it becomes difficult during cross validations. Having two evaluation metrics for a model becomes uneasy to make decision if we multiple models to evaluate.
-
- F1 Score : automatically supports for selecting optimum precision and recall.

F1 score

How to compare precision/recall numbers?

	Precision (P)	Recall (R)	Average	F ₁ score
Algorithm 1	0.5	0.4	0.45	0.444
Algorithm 2	0.7	0.1	0.4	0.175
Algorithm 3	0.02	1.0	0.501	0.0392

print("y=1")

~~Average = $\frac{P+R}{2}$~~

Harmonic mean

$$F_1 \text{ score} = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)} = 2 \frac{PR}{P+R}$$

▼ Powerful algorithms implementation ?

▼ Decision tree

- root node : Top most node on the tree
- leaf node : node other than root node
-

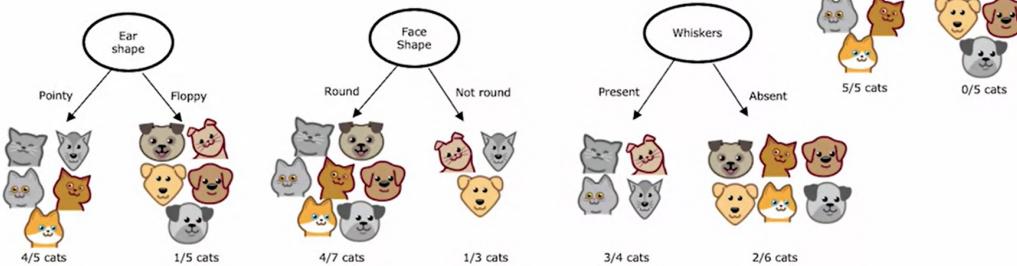
▼ How to select the root node, how is it calculated?

- The root node feature is selected based on maximizing purity or minimizing purity.
 - what is purity here?
 - If number of classes within the subdivided group is 1 then it is high purity else it leads to impurity.
 - For a high purity leaf, leaf will have only one class.
- The right top corner explains what purity is →

Decision Tree Learning

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

Maximize purity (or minimize impurity)



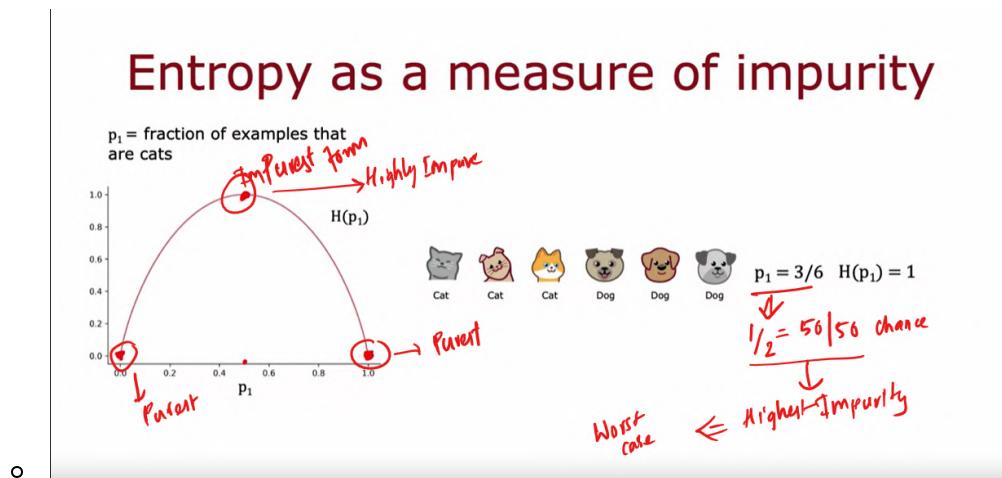
- How do we know when to stop the decision tree?
 - When purity is 100% meaning a leaf will have only one class
 - When tree exceeds the defined maximum depth it stops, as we know decision trees are greedy algorithm.

- When the improvements in purity of leaf is less than the defined purity on the code then the algorithm stops
- If the number of examples in any given node becomes less than the minimum defined example number then the algorithm stops.
-

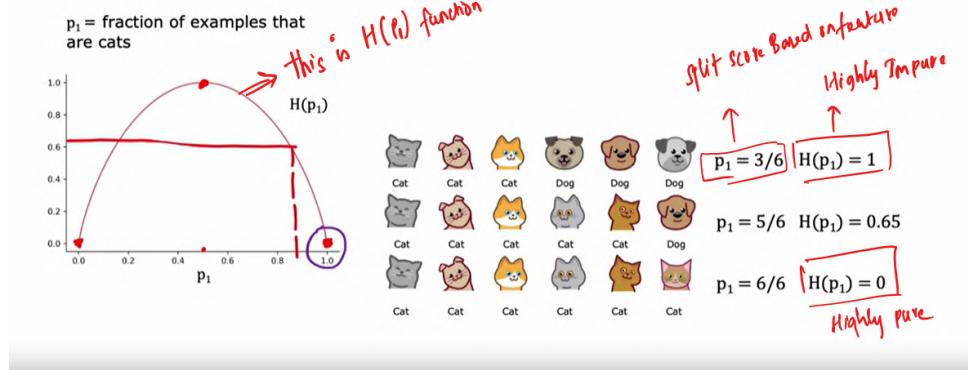
▼ How the split happens in the leaf node ?

▼ how to measure impurity?

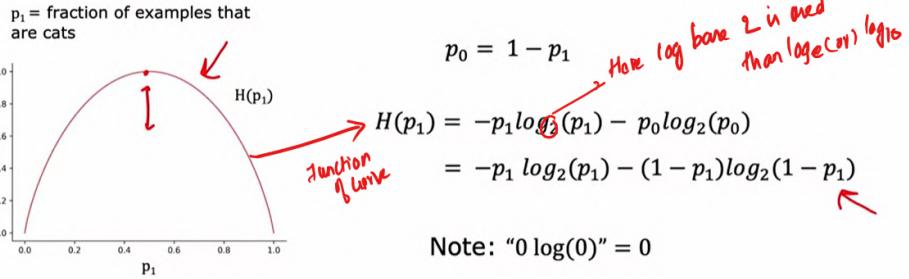
- entropy is a measure of impurity: if high entropy then high impurity, if low entropy then high purity
- Highest purity is when a leaf has only one class, highest impurity is when there are more than one class which are **equally distributed**
- what is entropy then?



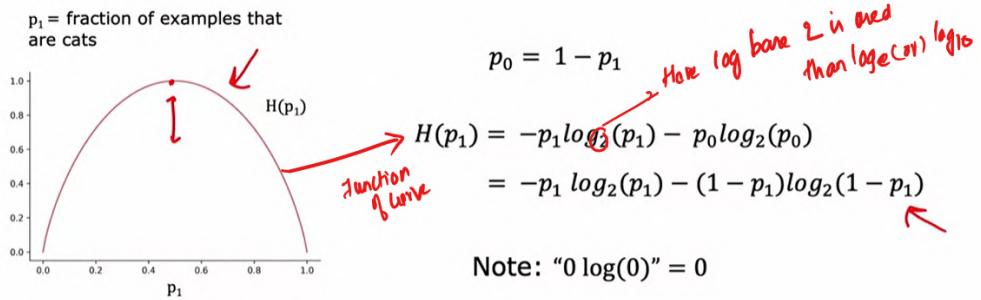
Entropy as a measure of impurity



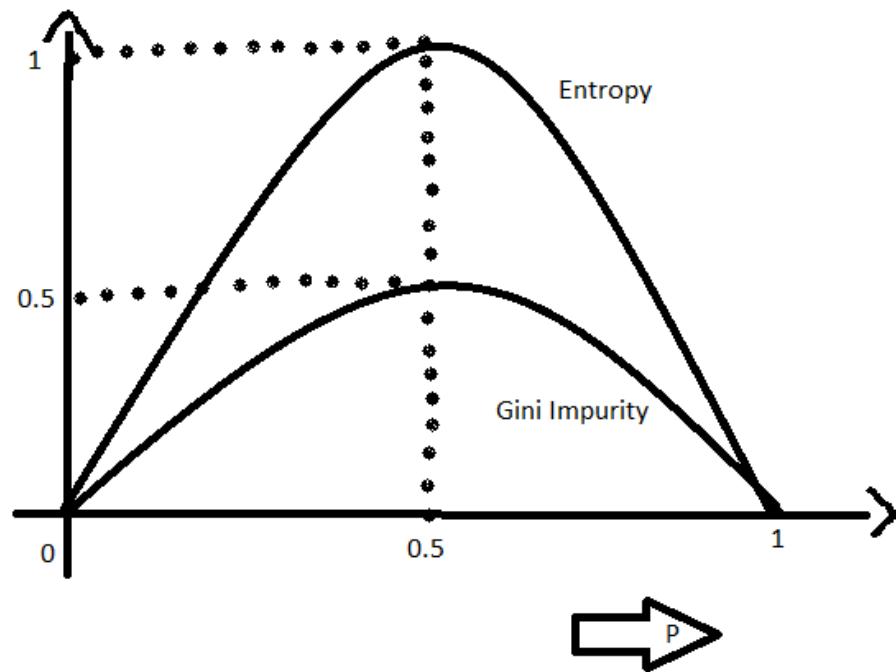
Entropy as a measure of impurity



Entropy as a measure of impurity



- Other functions for measuring impurity of decision tree :

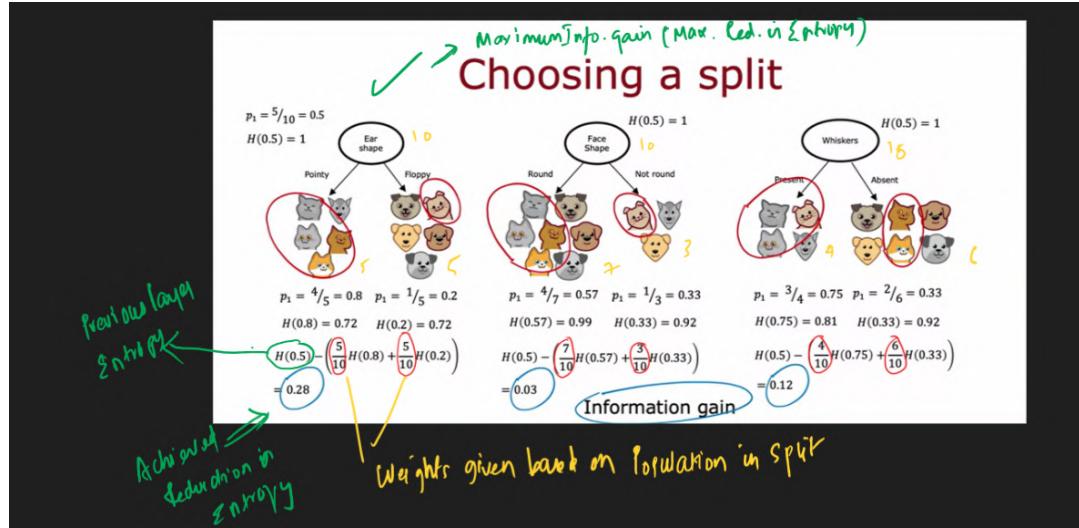


$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

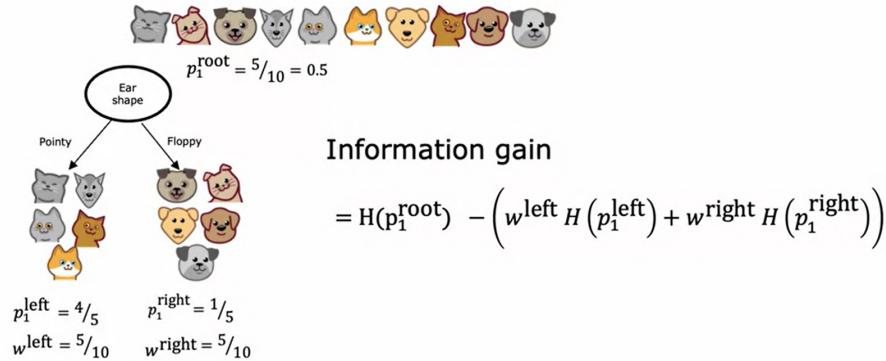
$$Gini(E) = 1 - \sum_{j=1}^c p_j^2$$

▼ How to choose the split based on information gain ?

- Information gain: The maximization or minimization of entropy (purity measuring function)
-



Information Gain



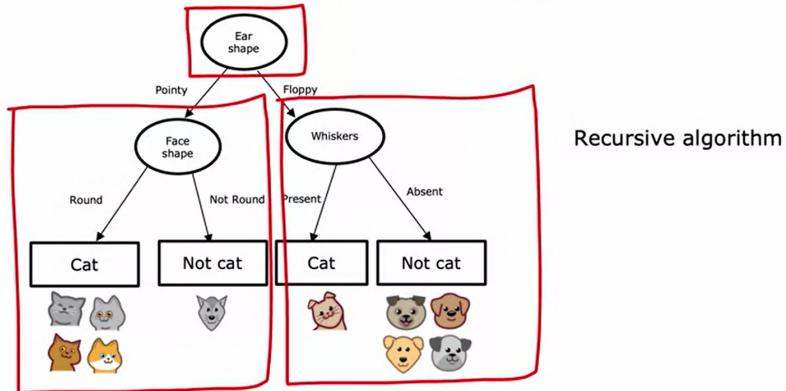
▼ How to build a decision tree ?

Decision Tree Learning

- Start with all examples at the root node
- Calculate information gain for all possible features, and pick the one with the highest information gain
- Split dataset according to selected feature, and create left and right branches of the tree
- Keep repeating splitting process until 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

- Recursive Algorithm: A specific algorithm is used from top to bottom of a tree, repetition of algorithm at scale.

Recursive splitting



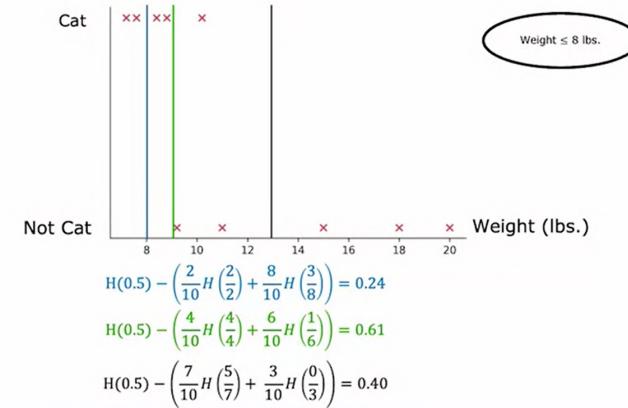
- How to select the depth of Decision tree
 - Select based on open source library
 - higher depths > larger tree > overfitting case might be.
 - Try validation techniques : for selecting the depth of tree, however open source libraries may be much better.
 - Use information gain to limit the growth of tree and number of examples within last leaf node.

▼ How to handle categorical features ?

- What will happen, if it has more than 2 categories in the same column → seems challenging for decision tree.
- So, we will handle by using one-hot encoding.

▼ how decision tree works for continuous range for values ?

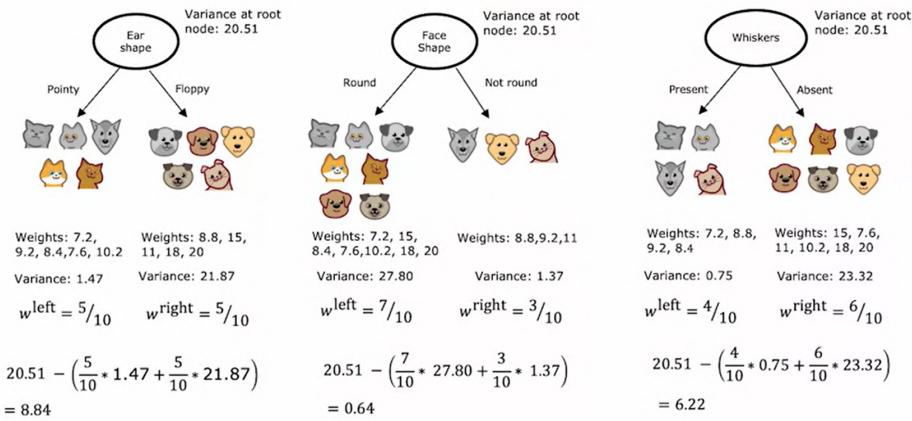
Splitting on a continuous variable



▼ how the generalization happens with decision tree ?

- we try to reduce the variance incase if it is a regression problem.
Since output is continuous variable , we try to reduce the variance in each level of trees:

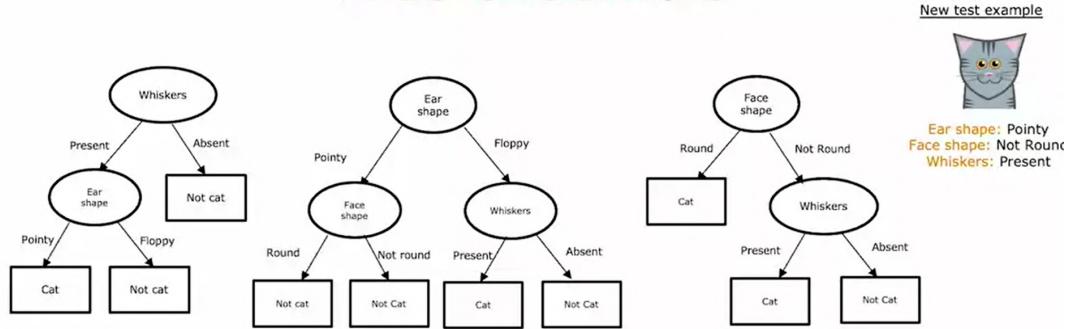
Choosing a split



▼ What is ensemble of decision tree

- combining multiple decision trees and consolidating the result.
-

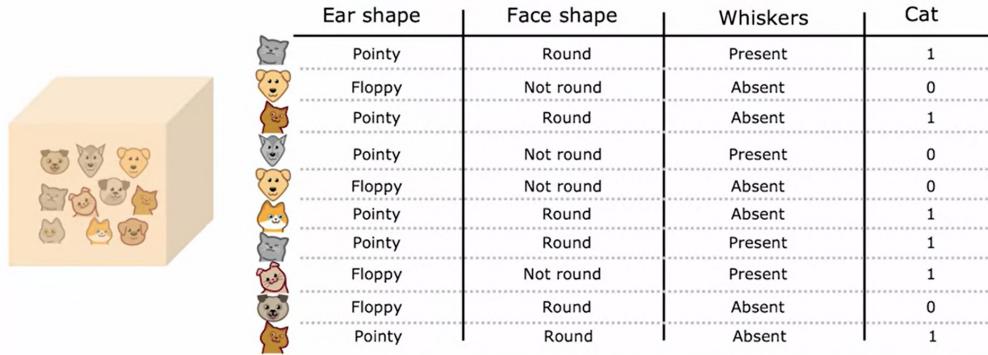
Tree ensemble



▼ sampling with replacement

- example means rows with all the features in.
- higher number of tree makes model generalized and the performance remains same.
- This is also called bagged decision trees.
- For improving the variety of decision trees to be made :
 - We can take sample with replacements (subset of rows)
 - we can randomizing the feature selection (subset of columns)
 - for higher number of features(columns) :
 - if n number of features are available in total, we select k number of feature to create decision tree.
 - if the number of features are too high then : $k = n^{(1/2)}$ → used generally. : this is called Random Forest algorithm.

Sampling with replacement



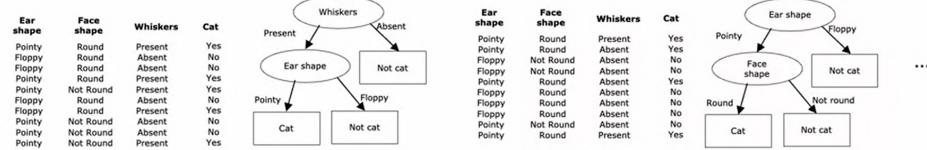
Generating a tree sample

Given training set of size m

For $b = 1$ to B :

Use sampling with replacement to create a new training set of size m

Train a decision tree on the new dataset



▼ How does the XGboost work ?

- Extreme gradient boosting
 - Has built in regularisation to reduce the problem of overfitting
 -

Boosted trees intuition

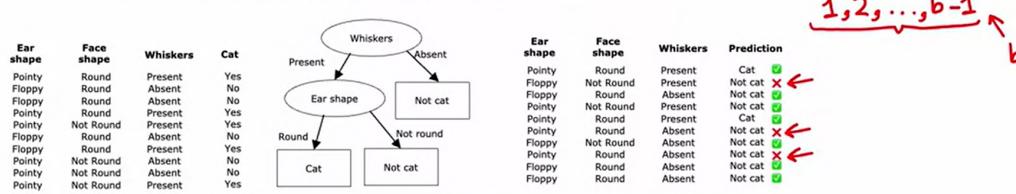
Given training set of size m

For $b = 1$ to B :

Use sampling with replacement to create a new training set of size m

But instead of picking from all examples with equal ($1/m$) probability, make it more likely to pick misclassified examples from previously trained trees

Train a decision tree on the new dataset



Using XGBoost

Classification

```
→from xgboost import XGBClassifier
→model = XGBClassifier()
→model.fit(X_train, y_train)
→y_pred = model.predict(X_test)
```

Regression

```
from xgboost import XGBRegressor
model = XGBRegressor()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

▼ Decision tree vs neural networks

-

Decision tree	Neural networks
Works best for structured data	Works well with all types of data and mixed data as well
Unsuited for unstructured data	Multiple models can be threaded together.
Faster training	May be slower based on the application

Smaller decision trees are interpretable, but does not work with ensembles.	Has transfer learning, which is quite efficient.
-----------------------------------------------------------------------------	--------------------------------------------------

