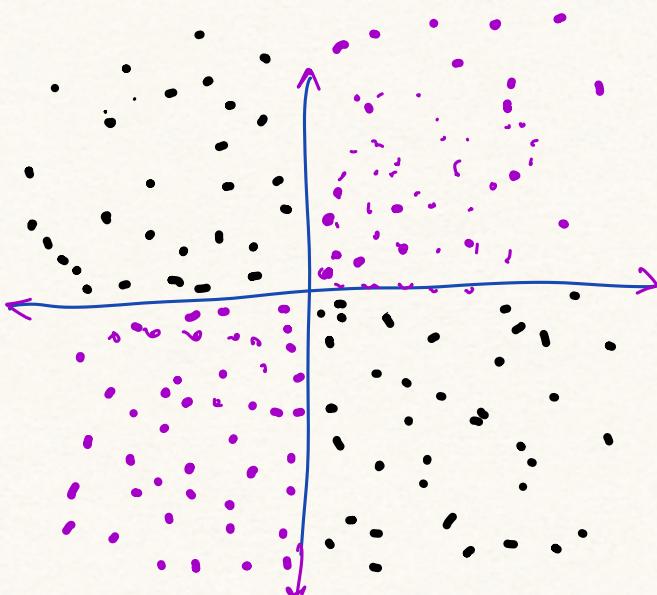


Neural Network



Non-linear Classification

problem because a linear function cannot separate all the blue dots from pink dots.

Hence Non-Linear means you cannot perfectly predict a label with a model of the form $b + w_1x_1 + w_2x_2$, means decision surface is not a line

if we perform feature cross on our features x_1 & x_2 we can then represent non-linear relationship between two features using linear model

$$b + w_1x_1 + w_2x_2 + w_3x_3$$

$$\text{where } x_3 = \underline{x_1 * x_2}$$

↳ feature cross between features x_1, x_2

\Rightarrow By adding feature cross x_1, x_2 , model can learn hypothesis shape that separates the blue dot & orange dots.

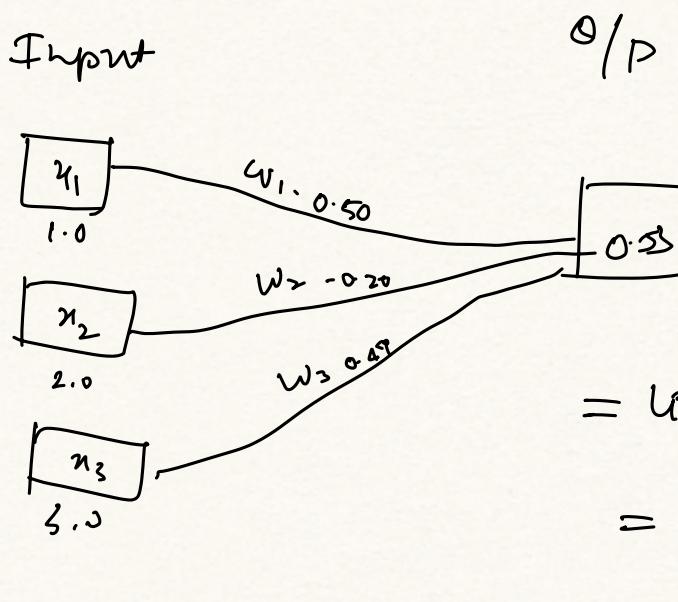
Hence \rightarrow Neural Network are family of model who

architecture designed to get non-linear patterns in data.
 During training Model get the optimal values for features used to
 minimize loss given input data

\Rightarrow Neural Networks learns non-linearity

$$\frac{y}{O/P} = b + w_1 x_1 + w_2 x_2 + w_3 x_3$$

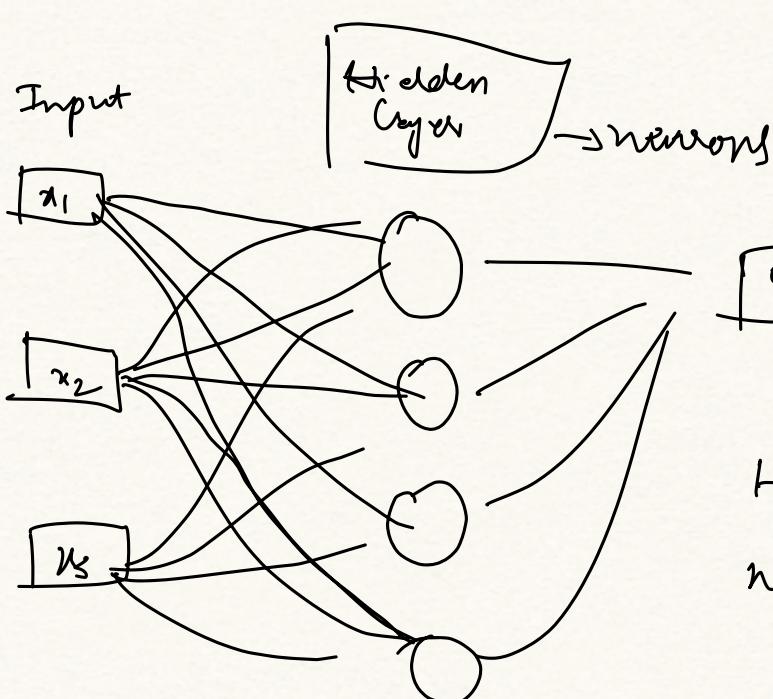
Input data



$$= \text{Linear}(w_1 x_1 + w_2 x_2 + w_3 x_3 + b)$$

$$= 1(0.50 + -0.41 + -1.46 + 0.83)$$

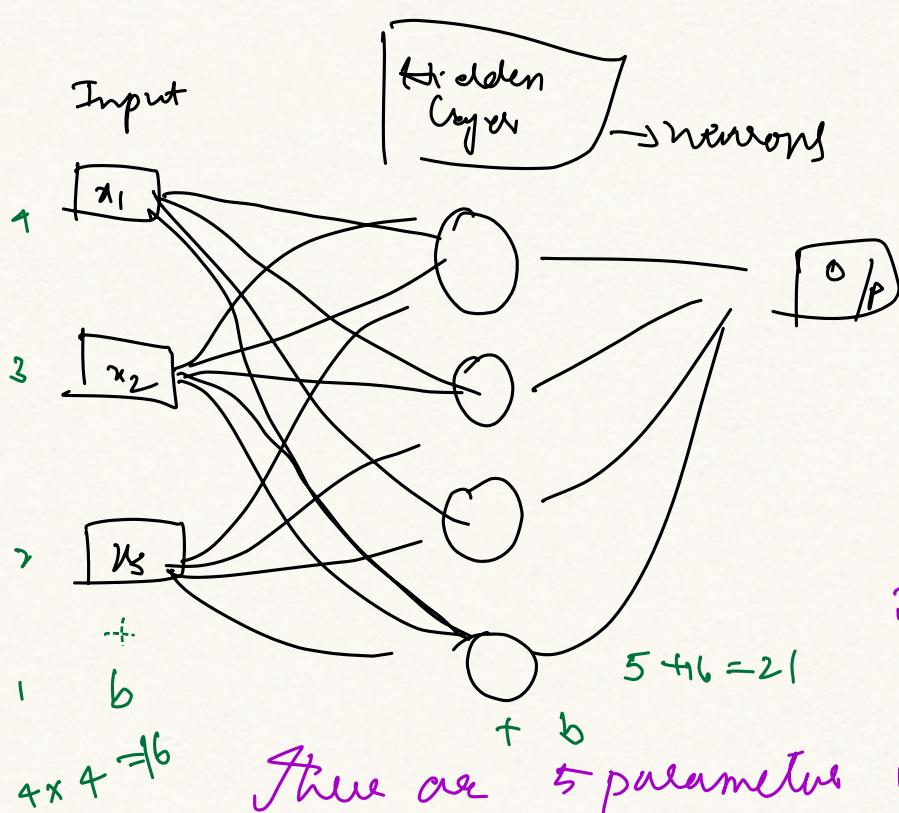
$$= \underline{\underline{0.53}}$$



Adding one more layer
 to network between Input & Output
 Hence This additional layer is called
 neurons

* The value of each neuron is calculated as the

Same way as the output of a linear Model, Sum of the product of each of its input (neurons of previous layer) \times unique weight parameter plus the bias, In the same manner neurons in next layer calculates the hidden layer's neuron value as input



How many parameters (weights and bias) does this neural network have

4 parameters used to calculate value of i in the hidden layer. i.e. 4 node values.

3 weights & bias

Sums = 16

4 weights & bias.

→ All this linear calculations i.e multiplication & addition will result a linear, Hence this model still doesn't calculate any non-linearities

Activation functions

To make linear network to learn non-linearity we have to introduce something like sigmoid function like we did in logistic regression

$$\text{Sigmoid} \left(\frac{1}{w_1 u_1 + w_2 u_2 + w_3 u_3} \right)$$

From here if we introduce Sigmoid function to get the value of the derivative

These sigmoid serves as activation functions for neural network, or nonlinear transform of a neuron's output value before the value is passed as input to calculate of the next layer of the neural network

Common activation functions

① Sigmoid

② tanh

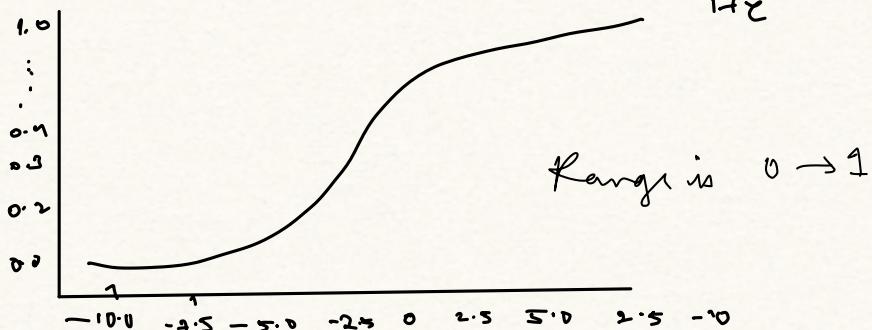
③ ReLU

* Sigmoid function → we use sigmoid in case

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

it is sigmoid as it is shaped, which squeezes value between 0-1

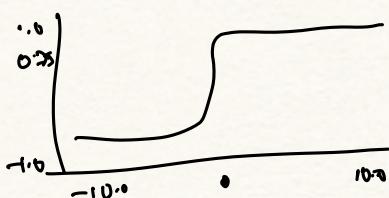
basically it is called as logistic function $F(x) = \frac{1}{1+e^{-x}}$



Range is $0 \rightarrow 1$

* The tanh (short for "hyperbolic Tangent") function transforms input x to produce an output value between $-1 < 1$

$$F(x) = \tanh(x)$$



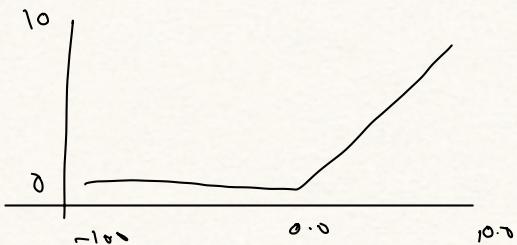
Range of tanh function is
-1 to 1

* ReLU - Rectified Linear Unit activation functions (or ReLU)

* if the input value x is less than 0, returns 0.

* value ' x ' is greater than or equal to 0, return the input

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



ReLU often works a little better as an activation function than smooth functions like sigmoid or tanh.

because it is less susceptible

to vanishing gradient problem during neural network training

Other functions

→ Mathematical functions can serve as an activation functions
→ $(w \cdot x + b)$

Keras provides out of box support for many activation functions
recommended is (ReLU)

Back Propagation

It is most common { algorithm } and { technique } used to train neural network, it makes
gradient descent feasible for multilayer neural network

Best Practices

Back propagation failure cases & most common ways to regularize
a neural network

Vanishing gradients

gradients in neural network layers can become very small, when
gradient values approach 0 for the lower layers, the gradient are said to vanish.
Layers with vanishing gradients train very slowly or not at all.

ReLU activation function can help

Exploding Gradients

If weight are very large, then the gradients for the layers involve

product of many large terms, gradient that get too large to converge

Batch normalization → can help exploding gradients

Dead ReLU Units

When weighted sum falls under zero, ReLU can stuck and contributes nothing to network's output, then gradients can no longer flow through it during back propagation

Learning Rate can solve or help ReLU units from dying

Dropout Regularization

→ another regularization is useful for neural networks,
it works randomly "dropping out" unit activations in a network
for single gradient step. The more you drop the more regularization

0.0 → No dropout regularization

1.0 → Drop out all nodes → Model learns nothing

Values between 0.0 and 1.0 → Most useful

Multiclass Classification

In Binary Classification

spam \leftrightarrow not spam

tumor is benign $\overset{\text{benign}}{\leftrightarrow}$ tumor is malignant or not malignant.

Multiclass

→识别狗 beagle, basset hound or blood hound

flower → Siberian Iris, Dutch Iris, Blue flag Iris or
Dwarf bearded Iris

plane → Boeing, Airbus etc

single \rightarrow apple, car, candy, dog etc.

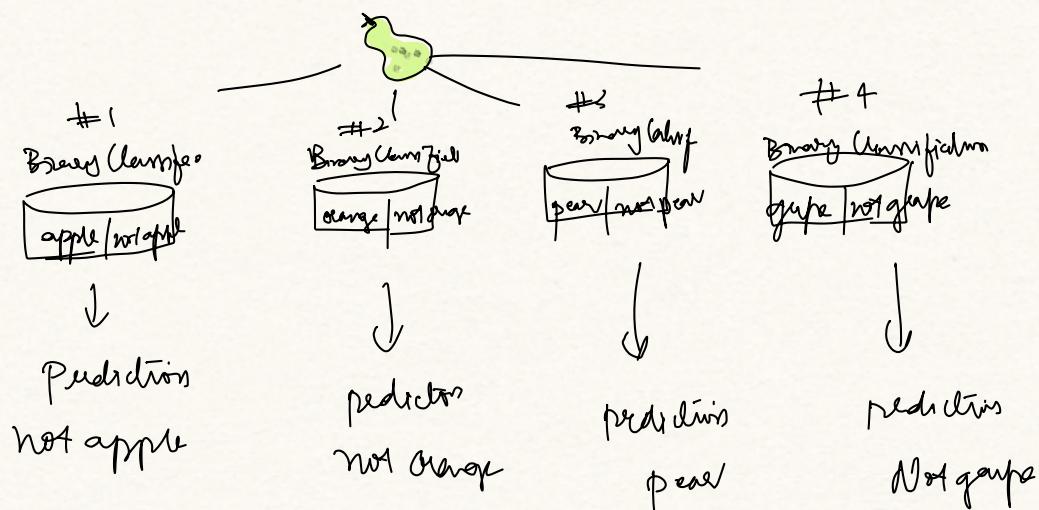
Variants of Multiclass

① One-vs-all

② One-vs-one \rightarrow also called as softmax.

① One-vs-all uses binary classification for a series of yes/no predictions across multiple possible labels

Eg given classification problem with N solutions a one-vs-all N-separate binary classifications \rightarrow one binary classifier for each problem or possible outcome



One Vs One (Soft Max)

In binary classification probability is compared are determined

independently from others.

\Rightarrow When we want to predict probability of every point to each other item relatively

is One Vs One

\Rightarrow For one vs one we apply function called softmax which assigns decimal probabilities to each class in multiclass problem such that all probabilities add up to 1.0.

$$e^{w_i^T x + b_i} / \sum e^{w_i^T x + b_i}$$

$$P(y = j/x) = \frac{e^{(\omega_j^T x + b_j)}}{\sum_{k \in K} e^{(\omega_k^T x + b_k)}}$$

it basically extends the formula $= \frac{1}{1 + e^{-y}}$

Softmax options

- * full softmax \rightarrow Softmax calculating probability for every possible class
 - * Condiate Sampling \rightarrow Softmax calculate probability for all positive labels
with only sample of negative labels
- Eg: if we want to determine weight or bound we don't need to give all the non-doggy stuff.

Softmax can be used to come class to one class and it expects to be similar class.