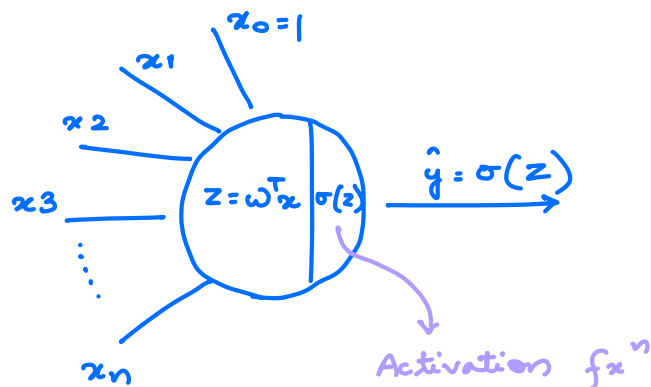


In this lecture we will see what is a Perceptron.

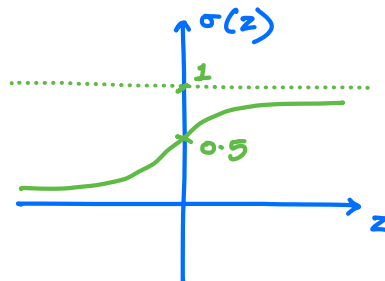
Perceptron is a single layer neural network. We will be learning about the loss function and later on we will derive weight update rule.



$$z: w^T x = [w^0 \ w^1 \ w^2 \ \dots \ w^n] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

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

↓
o/p is any no. b/w 0 and 1



let us say we are doing a classification problem in which we want to predict if given image is of dog or cat.

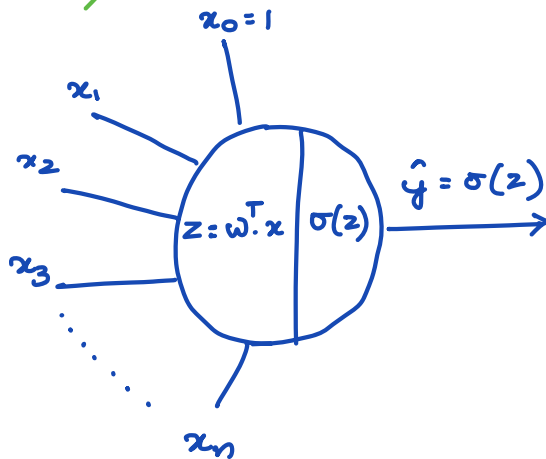


→ 0 or 1 ?

\hat{y} lies in range 0 and 1.

If $\hat{y} = (0.8) \approx 1$ (Round off to 1) } \hat{y} tells which class
 $\hat{y} = (0.3) \approx 0$ (Round off to 0) } a given input belongs to.

Perceptron will act as a linear classifier. It can do Binary Classification)



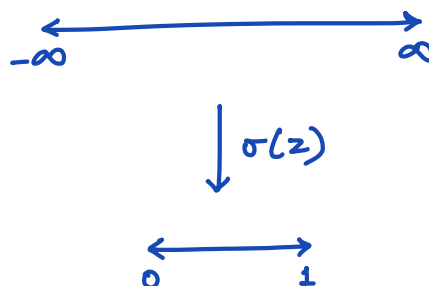
Basic Neural Network
with only 1 node
←

Goal of learning Algorithm is to learn parameters w , if you get data of type 1 then output should be close to 1 and if you get data of type 0 then output should be close to 0.

If input is very large, let's say $z=100$ then $\sigma(z)=1$

$z=-200$ then $\sigma(z)=0$

Activation $f(x)$ is compressing your number line in the range 0 to 1.



$\sigma(z)$ tells the probability a given input belongs to which class.

Let's say $\sigma(z) = 0.7$, you are 70% confident that it belongs to class 1, you are 30% confident that it belongs to class 0.

$$P(y=1) \rightarrow \sigma(z)$$

$$P(y=0) \rightarrow 1 - \sigma(z)$$

$$z = w^T \cdot x$$

This is how we can make predictions using a simple perceptron.

How to train a perceptron?

In every ML Algo, we do the following:

- Model (single node of Perceptron, when we combine multiple such nodes we get a neural network which is called as multi layer perceptron)

- Loss



Algo predicts probability, with what probability it belongs to class Dog.

$$P(y=1) : \hat{y} : 0.7 \quad 0.8 \quad 0.8 \quad 0.3 \quad 0.4$$

with 0.3 probability it belongs to class Dog, it belongs to class cat with probability 0.7

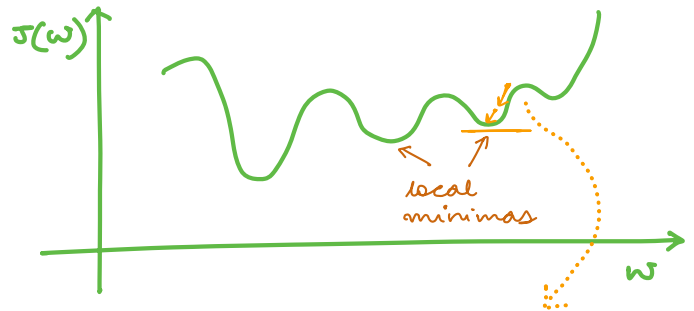
$$y : 1 \quad 1 \quad 1 \quad 0 \quad 0$$

$$J(w) \xrightarrow{\text{MSE}} \frac{1}{m} \sum_{i=1}^m \left(\hat{y}^{(i)} - y^{(i)} \right)^2$$

$$\hat{y}^{(i)} = \sigma(z) = \sigma(w^T x)$$

$\hat{y}^{(i)}$ is a function of w .

Problem with this kind of loss function is it is a non-convex function.



If you start from this point and use gradient descent then you will be stuck in local minima. There are multiple local minimas and you will stuck inside it.

To overcome this difficulty we use another kind of loss called as log loss.

log loss / Binary Cross Entropy

Proof for this is already covered in Maximum Likelihood Estimation for logistic Regression.

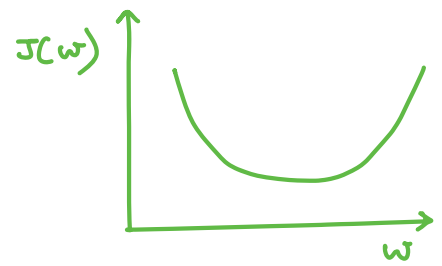
$$-\sum_{i=1}^m \left(y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)}) \right)$$

It is a convex $f(x)$

only 1 local and global minima

- if $y^{(i)} = 0$

then loss $f(x^n) = \log(1 - \hat{y}^{(i)})$



- if $y^{(i)} = 1$

then loss $f(x^n) = \log \hat{y}^{(i)}$

if predictions are not correct then we add log of misclassification.

- if $y^{(i)} = 0$ and $\hat{y}^{(i)} = 0$

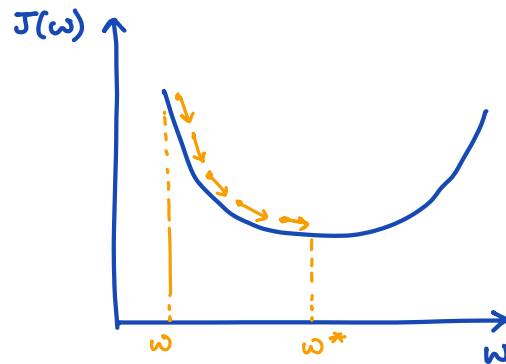
then loss $f(x^n) = 0$

→ that means when actual values matches with predicted value then you are adding 0 to loss.

- if $y^{(i)} = 1$ and $\hat{y}^{(i)} = 1$

then loss $f(x^n) = 0$

- Weight update Rule



Start with some random w and update w as:

$$w = w - \eta \left(\frac{\partial J}{\partial w} \right) \quad \left. \vphantom{\frac{\partial J}{\partial w}} \right\} \text{Gradient update Rule}$$

-ve
+ve

$$J(w) = - \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log (1 - \hat{y}^{(i)})$$

$$\hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(\omega^T x^{(i)})$$

$\sigma(\omega)$ is a fcn of $\hat{y}^{(i)}$

$\hat{y}^{(i)}$ is a fcn of $\omega^T x^{(i)}$