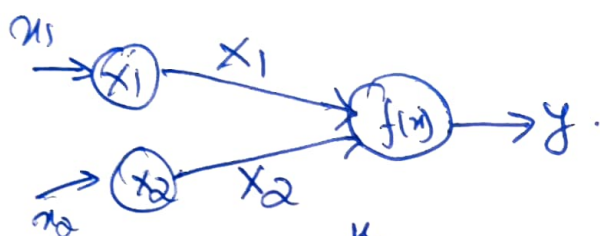


Problem 1

x_1	x_2	Class
1	1	+1
-1	-1	-1
0	0.5	-1
0.1	0.5	-1
0.2	0.2	+1
0.9	0.5	+1

$$W T x = 0$$

$$W_{\text{initial}} = [1, 1]$$



$$\text{Let } f(y(n)) \rightarrow \begin{cases} 1, & y(n) > 0 \\ 0, & y(n) = 0 \\ -1, & y(n) < 0 \end{cases}$$

$$y = b + \sum_{i=0}^n W_i x_i$$

① Let $b = 0$; learning rate, $\alpha = 1$.

8

$$y_{in} = W_1 x_1 + W_2 x_2$$

f.1. The Perceptron Learning Algorithm will Converge in 3 steps.

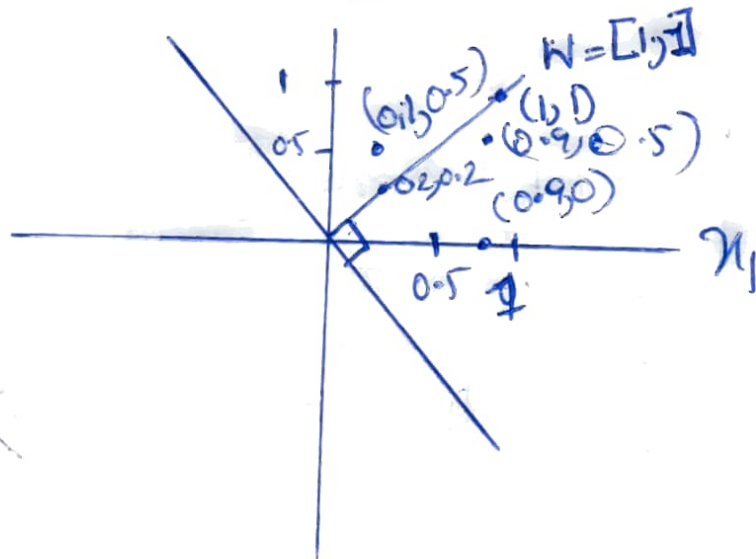
Step-1

Class	x_1	x_2	w	$W T x$	Correct	Corrected w	
1	1	1	[1, 1]	2	✓	—	$W T x = 2, f(x) > 1 \text{ (class)} > 0$
-1	-1	-1	[1, 1]	-2	✓	—	$W T x = -2, f(x) < 0, \text{ class} < 0$
-1	0	0.5	[1, 1]	0.5	X	$W = W_1 - \alpha$ $= [1, 0.5]$	$f(x) > 0, \text{ class} < 0$
-1	0.1	0.5	[1, 0.5]	0.35	X	$W = W_1 - \alpha$ $= [0.9, 0]$	$f(x) > 0, \text{ class} < 0$
1	0.2	0.2	[0.9, 0]	0.18	✓	—	$f(x) > 0, \text{ class} > 0$
1	0.9	0.5	[0.9, 0]	0.81	✓	—	$f(x) > 0, \text{ class} > 0$

$$W = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$W^T x = \begin{bmatrix} 1 & 1 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = x_1 + x_2 = 0$$

Group of equations, $L = x_1 + x_2 = 0$



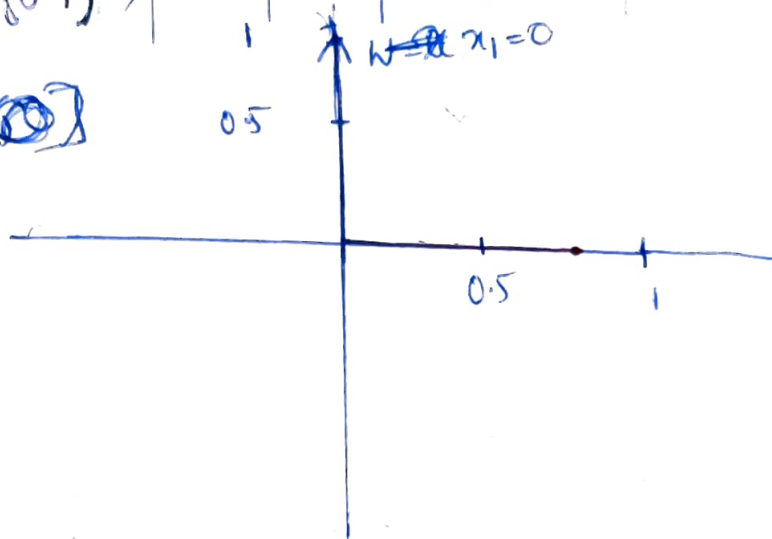
Class	x_1	x_2	W	$W^T x$	Correct	Corrected w
1	1	1	$(0.9, 0)$	0.9	✓	$f(x) > 0, \text{class}$
-1	-1	-1	$(0.9, 0)$	-0.9	✓	$f(x) > 0, \text{class}$
-1	0	0.5	$(0.9, 0)$	0	—	$[0.9, -0.5]$
-1	0.1	0.5	$(0.9, 0.5)$	-0.25	✓	$f(x) < 0, \text{class} < 0$
1	0.2	0.2	$(0.9, 0.5)$	0.08	✓	$f(x) > 0, \text{class} > 0$
1	0.9	0.5	$(0.9, 0.5)$	0.46	✓	$f(x) > 0, \text{class} > 0$

for $W = (0.9, 0)$

$$[0.9, 0] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$$

$$\Rightarrow 0.9x_1 = 0$$

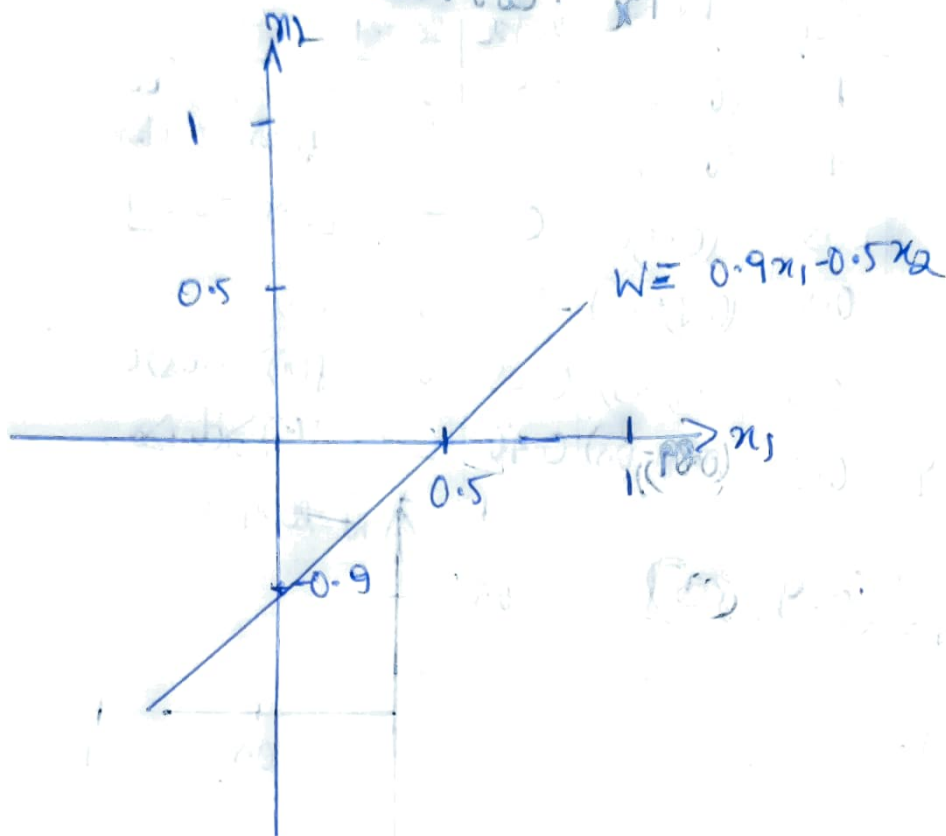
$$\Rightarrow x_1 = 0$$



Step 3

Class	x_1	x_2	W	$W^T x$	Correct	Correct w	
1	1	1	$(0.9, -0.5)$	0.4	✓	—	$f(x) > 0, \text{class} > 0$
-1	-1	-1	$(0.9, -0.5)$	-0.4	✓	—	$f(x) < 0, \text{class} < 0$
-1	0	0.5	$(0.9, -0.5)$	-0.25	✓	—	$f(x) < 0, \text{class} < 0$
-1	0.1	0.5	$(0.9, -0.5)$	-0.16	✓	—	$f(x) < 0, \text{class} < 0$
1	0.2	0.2	$(0.9, -0.5)$	0.08	✓	—	$f(x) > 0, \text{class} > 0$
1	0.9	0.5	$(0.9, -0.5)$	0.56	✓	—	$f(x) > 0, \text{class} > 0$

for $W = (0.9, -0.5)$



$$W = [0.9, -0.5]$$

$$W^T x = 0.9x_1 - 0.5x_2 = 0$$

$$\Rightarrow \underline{0.9x_1}$$

The final Decision boundary is $(0.9, -0.5)$.

③ Prove that Perceptron Learning Algorithm
Converges in a finite number of steps.

Soln

Definition:- Two sets P and N of points in an n -dimensional space are called absolutely linearly separable if $n+1$ real numbers w_0, w_1, \dots, w_n exist such that every point $(x_1, x_2, \dots, x_n) \in P$ satisfies $\sum_{i=1}^n w_i * x_i > w_0$ and every point $(x_1, x_2, \dots, x_n) \in N$ satisfies $\sum_{i=1}^n w_i * x_i < w_0$.

Proposition :- If the sets P and N are finite & linearly separable, the perceptron learning algorithm updates the weight vector w_T a finite number of times. In other words, if the vectors in P and N are tested cyclically one after the other, a weight w_f is found after a finite number of steps t which can separate the two sets.

Setup :-

- * If $x \in N$ then $-x \in P$. ($\because w^T x < 0 \Rightarrow w^T (-x) > 0$)
- * We can thus consider a single set $P' = P \cup N$ and for every element $p \in P'$ ensure that $w^T p > 0$.
- * Let us will normalize all the p 's so that $\|p\| = 1$.
- * Let w^* be the normalized solution vector.

Proofs-

* Now suppose at time step t we inspected the point P_i , and found that $W^T \cdot P_i \leq 0$.

* we make a correction $W_{t+1} = W_t + P_i$

* let β be the angle between W^* and W_{t+1}

$$\cos \beta = \frac{\cancel{W^* \cdot W_t} + W^* \cdot W_{t+1}}{\|W_{t+1}\|}$$

$$\begin{aligned} \text{Numerator} &= W^* \cdot W_{t+1} = W^* (W_t + P_i) \\ &= W^* W_t + W^* P_i \end{aligned}$$

$$\geq W^* W_t + \delta \quad (\delta = \min \{W^* \cdot P_i \mid P_i \in \mathcal{V}\})$$

$$\geq W^* (W_{t-1} + P_j) + \delta$$

$$\geq W^* W_{t-1} + W^* P_j + \delta$$

$$\geq W^* W_{t-1} + 2\delta$$

Every time we make a correction a quantity δ is added. Say if K corrections are made, a quantity of $K\delta$ is added.

$$\cos \beta = \frac{W^* \cdot W_{t+1}}{\|W_{t+1}\|} \quad (\text{definition})$$

$$\begin{aligned} \text{Numerator} &\geq W^* W_0 + K\delta \quad (\text{proved by induction}) \\ \text{Denominator} &= \|W_{t+1}\|^2 \end{aligned}$$

$$\begin{aligned} &= (W_t + P_i)^2 \\ &= \|W_t\|^2 + 2W_t P_i + \|P_i\|^2 \leq \|W_t\|^2 + \|P_i\|^2 \quad (\because W_t P_i \leq 0) \end{aligned}$$

$$\leq (\|w_{t-1}\|^2 + 1)$$

$$\leq \|w_{t-1}\|^2 + 1 + 1$$

$$\text{Denominator} \leq \|w_{t-1}\|^2 + (K)$$

$$\cos \beta \geq \frac{w^* w_0 + K}{\sqrt{\|w_0\|^2 + K}}$$

$\cos \beta$ thus proportional to \sqrt{K}

i.e. K has to be finite so that $\cos \beta \leq 1$.

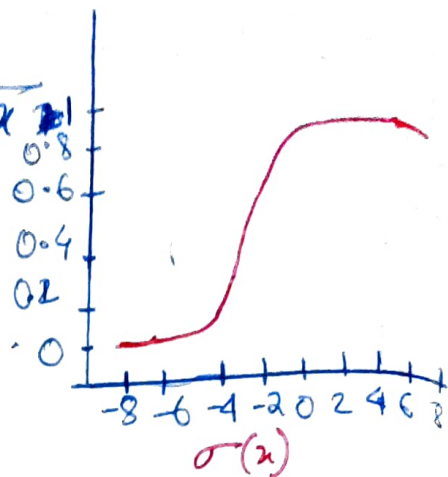
∴ Thus there can only be a finite number of corrections (K) to w and the perceptron algo will converge.

③ Compute Derivative of following activation functions:-

① Sigmoid.

$$\text{Sigmoid}(x) = \sigma(x) = \frac{1}{1+e^{-x}}$$

$$\frac{d}{dx} \sigma(x) = \frac{d}{dx} \frac{1}{1+e^{-x}}$$



Applying Quotient rule:-

$$\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{\left(g(x) \frac{d}{dx} f(x) \right) - \left(f(x) \frac{d}{dx} g(x) \right)}{\{g(x)\}^2}$$

$$\frac{d}{dx} \sigma(x) = \frac{(1+e^{-x}) \frac{d}{dx} 1 - 1 \times \frac{d}{dx} e^{-x}}{(1+e^{-x})^2}$$

$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \frac{1+e^{-x} - 1}{(1+e^{-x})^2}$$

$$= \frac{1+e^{-x}}{(1+e^{-x})^2} - \frac{1}{(1+e^{-x})^2}$$

$$= \frac{1}{1+e^{-x}} - \frac{1}{(1+e^{-x})^2}$$

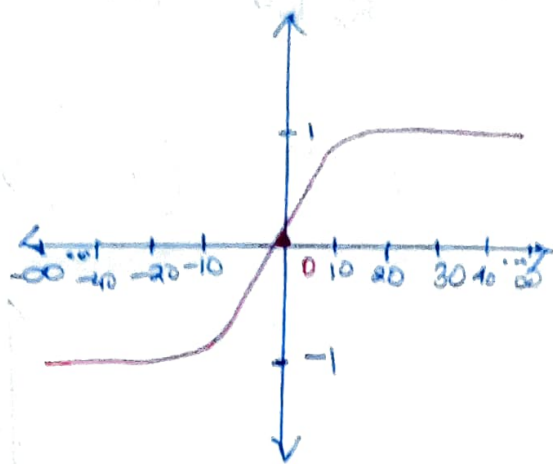
$$= \frac{1}{1+e^{-x}} \left\{ 1 - \frac{1}{1+e^{-x}} \right\}$$

$$= \sigma(x)(1 - \sigma(x))$$

$$\therefore \frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x)) //$$

(11) \tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \\ = T(x)$$



$$f(x) = \tanh(x)$$

$$\frac{d}{dx} T(x) = \frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Applying quotient rule

$$\frac{d}{dx} \frac{f(x)}{g(x)} = \frac{g(x) \frac{d}{dx} f(x) - f(x) \frac{d}{dx} g(x)}{g(x)^2}$$

$$\frac{d}{dx} \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^x + e^{-x} \left(\frac{d}{dx} (e^x - e^{-x}) \right) - (e^x - e^{-x}) \left(\frac{d}{dx} (e^x + e^{-x}) \right)}{(e^x + e^{-x})^2}$$

$$= \frac{(e^x + e^{-x})(e^x + e^{-x}) - \{(e^x - e^{-x})(e^x - e^{-x})\}}{(e^x + e^{-x})^2}$$

$$= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

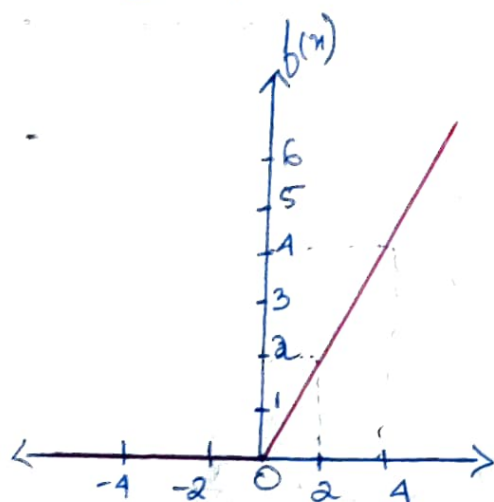
$$= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

$$\frac{d}{dx} T(x) = 1 - (T(x))^2$$

$$\therefore \frac{d}{dx} \tanh = 1 - (\tanh)^2$$

14) $R_e L U = \max(0, x)$

$$R_e L U = \begin{cases} x, & x > 0 \\ 0, & x \leq 0 \end{cases}$$



$f(x) = R_e L U$

~~At $x > 0$~~

$$f'(x) = \lim_{x \rightarrow 0} \frac{\max(0, x + \Delta x) - \max(0, x)}{\Delta x}$$

At $x > 0$

$$f'(x) = \lim_{x \rightarrow 0} \frac{x + \Delta x - x}{\Delta x} = 1$$

At $x < 0$

$$f'(x) = \lim_{x \rightarrow 0} \frac{0 - 0}{\Delta x} = 0$$

At $x = 0$

$$f'(0) = \lim_{x \rightarrow 0^+} \frac{0 + \Delta x - 0}{\Delta x} = 1$$

$$f'(0) = \lim_{x \rightarrow 0^-} \frac{0 - 0}{\Delta x} = 0$$

LHL \neq RHL.
Left hand limit \neq Right hand limit.

$$\therefore \frac{d}{dx} \text{ReLU} = \frac{d}{dx} \max(0, x)$$

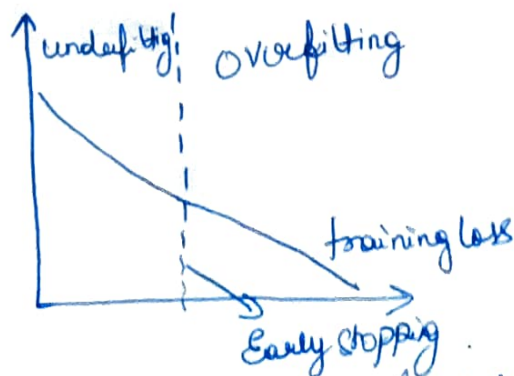
$$= \begin{cases} 1, & x > 0 \\ 0, & x < 0 \\ \text{Does not exist}, & x = 0 \end{cases}$$

② What are the strategies you will follow to avoid over fitting in a neural network.

Ans Overfitting is a situation where the model of neural network leaves the training data very well but does not perform well on a testing dataset.

① Simplifying the model by reducing the number of internal layers and making the network smaller. Accordingly change the input or reduction of internal layers.

② Early stopping:- Stopping the learning or training process of the model before it ~~has~~ crosses over-fitting.



However it is important to find the right point to stop to avoid underfitting or overfitting.

③ Data Augmentation :- By making small changes in the existing dataset, we make the model consider each augmented model as a distinct data set. This will make the neural network avoid overfitting.

④ Regularization :- It can be used to reduce the complexity of a model. The most common way is to add penalty to the loss function in proportion to the weight of the model.

$$\textcircled{3} \quad f(z) = z_1 x + z_2 y, \quad x = [1, 1]^T, \quad y = [1, 1]^T.$$

$$f: \mathbb{R}^2 \Rightarrow \mathbb{R}^2 \quad f(z) = z_1 x + z_2 y$$

$$z = g(r) = [r^2, r^3]$$

$$\frac{\partial f}{\partial r} = \frac{\partial f}{\partial z} \times \frac{\partial z}{\partial r} = \begin{bmatrix} x \\ y \end{bmatrix} [2r^2, 3r^2]$$

$$= \begin{bmatrix} 2r^2 x_1 & 3r^2 y_1 \\ 2r^2 x_2 & 3r^2 y_2 \end{bmatrix} \bigg|_{\substack{x=[1,1] \\ y=[1,1] \\ r=2}} = \begin{bmatrix} 8 & 12 \\ 8 & 12 \end{bmatrix}$$

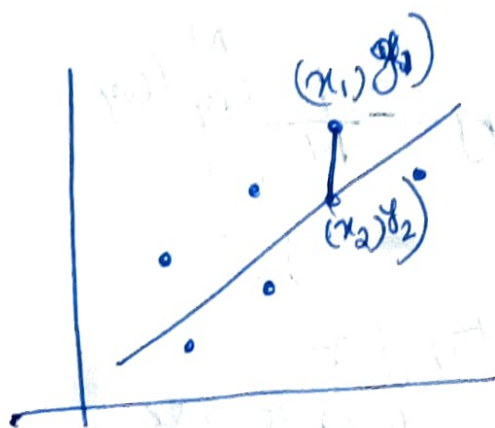
④ Define following loss functions in brief.

① Mean Squared Error (MSE)

~~A loss~~

Mean Squared Error (MSE) calculates the mean of distance of set data points from a regression line.

It is basically the average of errors squared from data.



$$\text{Error} = y_1 - y_2$$

$$\text{Squared Error} = (y_1 - y_2)^2$$

$$\text{Mean Squared Error} = \frac{1}{n} (y - \hat{y})^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

② Binary Cross-Entropy (BCE)

It compares each predicted probabilities to actual class output which can either be 0 or 1. It then calculates the score that penalizes the probabilities based on the distance from the expected value. Binary entropy is the negative

of the log of corrected predicted probabilities.

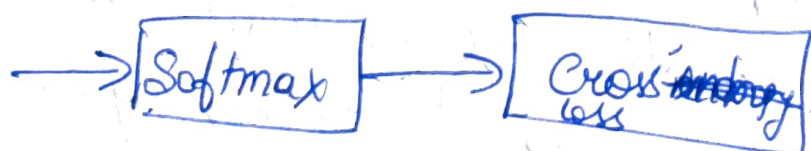
Eg	D	Actual	Predicted prob	Corrected
1	1	1	P1	P1
2	1	1	P2	P2
3	1	1	P3	P3
4	0	0	P4	P4
5	0	0	P5	P5

$$\text{Binary Cross Entropy} = -\frac{1}{N} \sum_{i=1}^N \log(p_i)$$

③ Categorical Cross Entropy :-

It is a loss function that is used in multi class classification tasks.

These are tasks where an example can only belong to one out of many possible categories and the model must decide which one.



$$f(s) / \text{Softmax} = \frac{e^{s_i}}{\sum_j e^{s_j}}$$

$$\text{Cross Entropy} = - \sum_{i=1} t_i \log(f(s)_i)$$

$$= - \log \left(\frac{e^{s_p}}{\sum_j e^{s_j}} \right)$$

Eg:- True label: Rabbit

Prediction: Dog=1, Cat=1, Rabbit=8,

Squirrel=2.

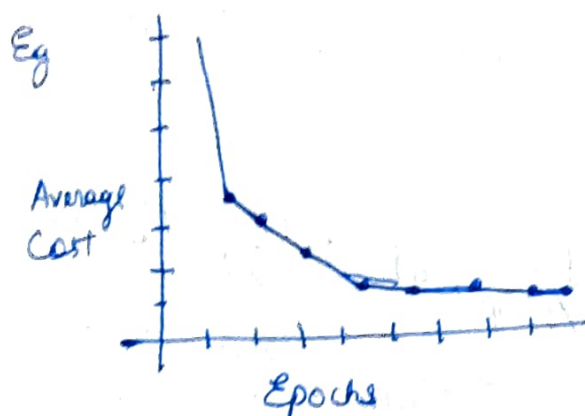
$$\text{Softmax} = \text{Dog} = e^1 / \text{sum}, \text{Cat} = e^1 / \text{sum}, \text{Rabbit} = e^8 / \text{sum}$$

$$\text{Squirrel} = e^2 / \text{sum} \quad // \text{sum} = e^1 + e^1 + e^8 + e^2$$

$$\begin{aligned} \text{Cross Entropy} &= - (0 \times \ln(\text{Dog}) + 0 \times \ln(\text{Cat}) \\ &\quad + 1 \times \ln(\text{Rabbit}) + 0 \times \ln(\text{Squirrel})) \\ &= \ln(\text{Rabbit}) \end{aligned}$$

⑤ ① Batch Gradient Descent:

Gradient descent is an optimization technique for minimising a function. In Batch Gradient Descent, all the training examples are used to calculate the mean gradient to update the parameters. It gives the best ~~result~~ ^{result} in case of convex or relatively smooth function.



The graph of Cost vs Epochs is also quite smooth because we are averaging over all the gradients of training data for a single step.

⑥ ② Stochastic Gradient Descent

Stochastic gradient descent attempts to find the global minimum by adjusting the configuration of the network after each training point. Instead of decreasing the error, or finding the gradient, for the entire data set, this method merely decreases the error by approximating the gradient for a randomly selected

It can adjust the network parameters in such a way as to move the model out of a local minimum and towards a global minimum.

③ Mini Batch Gradient Descent:-

Mini Batch Gradient Descent is a variation of the gradient descent algorithm that splits the training dataset into small batches that are used to ~~collect~~ calculate model error and update model coefficients.

Implementation may choose to sum the gradient over the mini batch which further reduces the variance of the gradient.

Mini Batch gradient descent seeks to find a balance between the noisiness of Stochastic gradient descent and efficiency of batch gradient descent. It is the most common implementation of gradient descent ~~used~~ in the field of deep learning.

④ Momentum based Gradient Descent:-

An adaptive Optimization algorithm uses exponentially weighted averages of gradients over previous iterations to stabilize the convergence, ~~res~~ and dealing with noisy data. To avoid the effect of noise and reduce its interference the data are fed in batches during optimization. This problem can be tackled using Exponentially

Weighted Averaged. $\gamma_0 = \beta \gamma_0 + (1-\beta) \gamma_t$.

This method is used as strategy in momentum based gradient descent to make it robust against noise in data samples.

⑤ Adam Solver :

Adam Solver optimization algorithm is an extension to Stochastic gradient descent that recently seen broader adoption for deep learning applications.

The learning rate is maintained for each network weight and separately adapted as learning unfolds.

It used both ~~the~~ Adaptive Gradient algorithm and Root mean square propagation.

~~It is also make~~.

The algorithm calculate an exponential moving average of the gradient and the squared gradient, and the parameters control the decay rate of these their moving averages.