

Q3

x_0	x_1	x_2	z
1	0	0	1
1	1	1	1
1	1	0	0

let $w_0 = 1$ $w_1 = -1$ $w_2 = 1$

let $x_0 = (2x_1 - 1)^2$

This problem is similar to XOR problem so the points are not linearly separable, so we have taken x_0 as quadratic function.

$$w_0 x_0 + w_1 x_1 + w_2 x_2 \text{ for } (100) = 1$$

$$w_0 x_0 + w_1 x_1 + w_2 x_2 \text{ for } (111) = 1$$

$$w_0 x_0 + w_1 x_1 + w_2 x_2 \text{ for } (110) = 0$$

If we make all three weights positive then in all cases we get positive output.

If we make all three negative then we get negative output in all cases.

If we make first two positive i.e. w_0, w_1 and last one negative we incorrectly classify third case. (110)

If we make first two negative i.e. w_0, w_1 and last one positive we incorrectly classify third case.

Q5

1. The filter size is $2 \times 2 \times 4$ for each filter.
There are in total 3 filters so no of channels = 12.

2.

0	1	0	1	0	1	0
0	0	1	0	1	0	0
0	1	0	1	0	1	0
0	0	1	0	1	0	0
0	1	0	1	0	1	0

output:-

0	0	0
0	0	0

3.

0	0
0	0

 (trivial solution)

1	0
0	-1

 (non trivial solution)

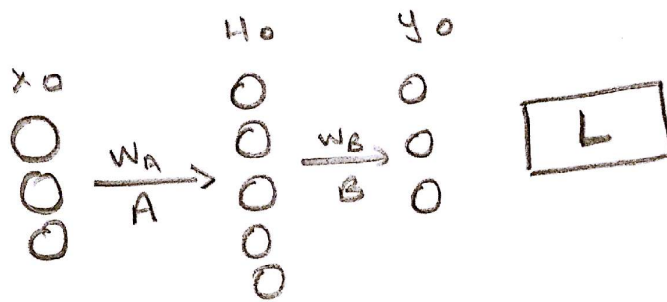
4.

1	-1
0	0

5. 4 filters were applied on layer 2 to get the feature maps in layer 3.

6. The first two feature maps in layer 3, I guess let us know that there might be a straight line at centre and a diagonal from left corner in image of layer 2.

Q2



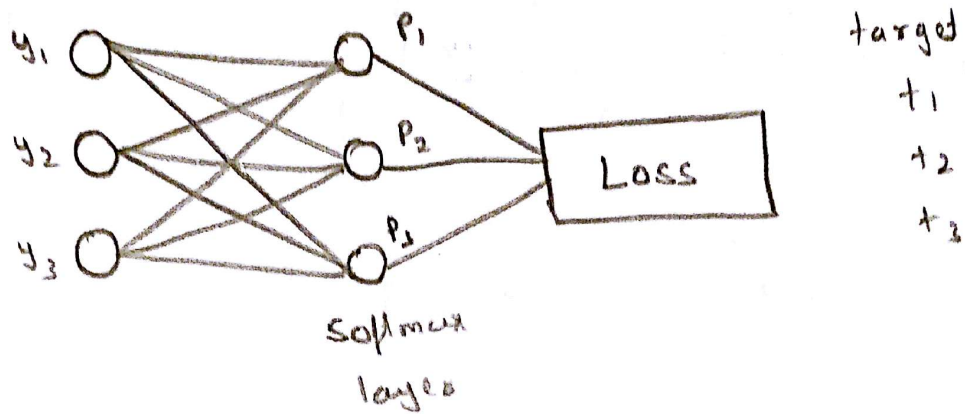
$$1) J_x = J_B \times W_B^T \times J_A \times W_A^T$$

$$2) \text{ Dimension of } J_x = [\text{len}(y_0), \text{len}(x_0)]$$

$$3) E_y = \frac{\partial L}{\partial y_0}$$

$$\frac{\partial L}{\partial x_0} = E_y \times J_x^T$$

option-2



Cross Entropy Loss $L = - \sum_i t_i \log(p_i)$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial p} \times \frac{\partial p}{\partial y}$$

$$\frac{\partial L}{\partial p} = - \sum_i t_i \times \frac{1}{p_i}$$

Now $p_i = \frac{\exp(y_i)}{\sum_{k=1}^N \exp(y_k)}$

$$\frac{\partial p_i}{\partial y} = \frac{\frac{\partial \exp(y_i)}{\sum_{k=1}^N \exp(y_k)}}{\frac{\partial y}{\partial y}}$$

Now if $i = j$ $\frac{\partial \exp(y_i)}{\partial y(i)} = \exp(y_i)$

if $i \neq j$ $\frac{\partial \exp(y_i)}{\partial y(j)} = 0$

For $i = j$ $\frac{\partial p}{\partial y} = p_i (1 - p_j)$

$i \neq j$ $\frac{\partial p}{\partial y} = -p_i \cdot p_j$

$$\frac{\partial L}{\partial y_i} = \left(- \sum_k t_k \times \frac{1}{p_k} \right) \times \left(p_i (1 - p_k) + (-p_i p_k) \right)$$

\swarrow when $k=i$ \nwarrow $k \neq i$

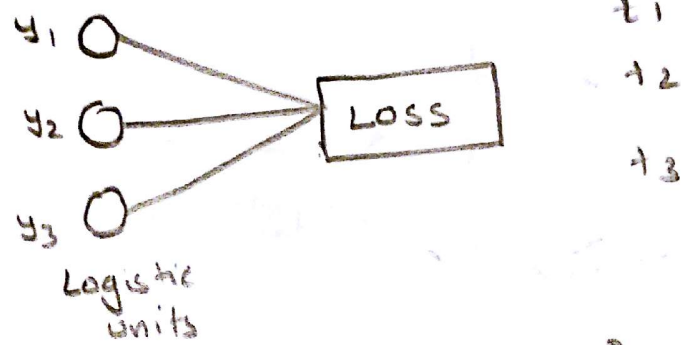
$$= -t_i (1 - p_i) + \sum_{k \neq i} t_k \times p_i$$

$$= -t_i + p_i t_i + \sum_{k \neq i} t_k \times p_i$$

$$\frac{\partial L}{\partial y_i} = p_i - t_i \quad \left(\sum_{k \neq i} t_k + t_i \right) = 1 \text{ because}$$

+ is one not vector

Option - 3



$$Loss = \frac{1}{2} \sum_{i=1}^n (y_i - t_i)^2$$

$$y = \sigma(z)$$

$$y = \frac{1}{1 + e^{-z}}$$

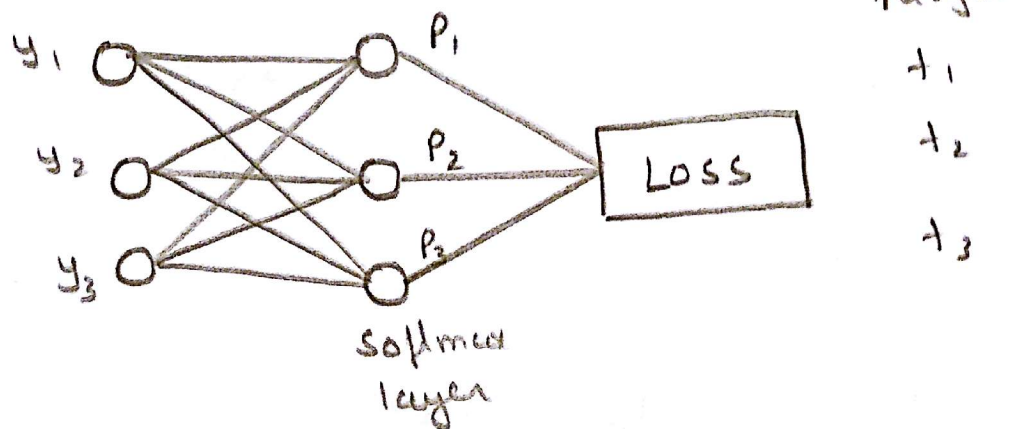
$$\frac{\partial Loss}{\partial z} = \frac{\partial Loss}{\partial y} \times \frac{\partial y}{\partial z}$$

$$\frac{\partial Loss}{\partial y} = \sum_{i=1}^n (y_i - t_i)$$

$$\frac{\partial y}{\partial z} = y(1 - y)$$

$$\frac{\partial Loss}{\partial z} = (y - t) \times y(1 - y)$$

option -2



$$\text{Loss} = \frac{1}{2} \sum_{i=1}^n (p_i - t_i)^2$$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial p} \times \frac{\partial p}{\partial y}$$

$$\frac{\partial L}{\partial p} = \sum_{i=1}^n (p_i - t_i)$$

$$p = \frac{\exp(y_i)}{\sum_{k=1}^n \exp(y_k)}$$

$$\begin{aligned} \frac{\partial p}{\partial y} &= p_i (1 - p_j) & i = j \\ &= -p_i \times p_j & i \neq j \end{aligned}$$

$$\begin{aligned} \frac{\partial \text{Loss}}{\partial y} &= \sum_k p_k - t_k \times (p_i (1 - p_k) - p_i p_k) \\ &= (p_i - t_i) \times p_i (1 - p_i) + \sum_{k \neq i} (p_k - t_k) \times (-p_i p_k) \end{aligned}$$

Identity, MSE

From the notes given in class it shows an example where our model confidently predict a negative label with $z = -5$, which gives $y = 0.0067$, for a positive example ($t = 1$). Plugging these values we get $\frac{\partial L}{\partial z} = -0.0066$. This is pretty small value in comparison to our mistake. It shows that the more confident the wrong prediction the smaller the gradient is.

Softmax, Cross Entropy

Here we get the output $\frac{\partial L}{\partial y}$ similar to the linear regression. So if $y > t$ you made too positive a prediction, so you want to shift prediction in negative direction. Conversely if $y < t$, you want to shift your prediction in positive direction.

Softmax, MSE

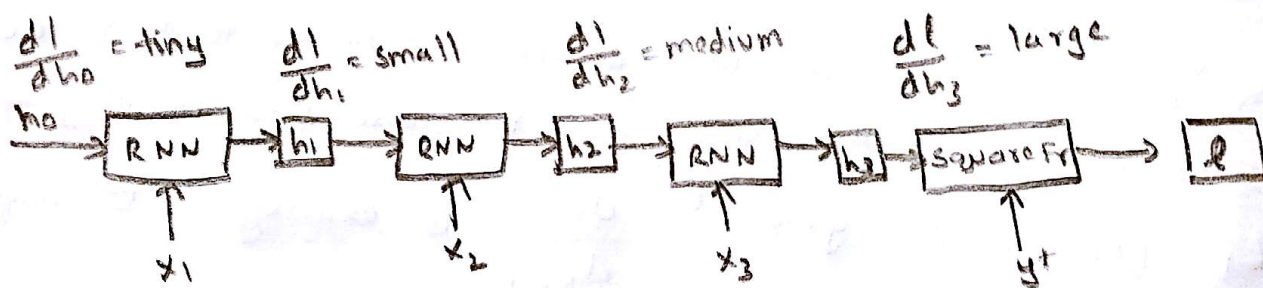
The problem with MSE in the classification task is that does not distinguish bad predictions from extremely bad predictions. If $t = 1$ then prediction of $y = 0.01$ has roughly same loss as for prediction of $y = 0.00001$, even though the latter is more wrong.

So I would suggest Jim to use softmax cross entropy loss.

Q4

Here there is verb "go" and with that there is dependency "library". Then people borrow books from library "so you have borrowed all of shakespeare's work". Then you have verb "read" which is also dependent on "library" as it is place to read. Now at last we have "bench" which is dependent on verb "read".

The sentence starts with "I" and in the last sentence we see that dependency "I will read....".



RNN makes it possible to model long distance dependencies because they have ability to pass information between timesteps.

Ex:- If some node h_{t+1} encode the information that "the subject of sentence is male", it is possible to pass this info to h_t which in turn can pass to h_{t+1} .

Suppose we have a RNN that predicts after several time steps and calculates loss that is expected to back propagate over all time steps. At each time step we run back propagation, the gradient gets smaller and smaller, and by the time we get back to beginning time steps, we have a gradient that is so small that it has no significant effect on parameters to be updated.

The reason this happens because we have either $0 < \frac{dh}{dh_{t-1}} < 1$ or $\frac{dh}{dh_{t-1}} > 1$ and we want $\frac{dh}{dh_{t-1}} = 1$. We do not have $\frac{dh}{dh_{t-1}}$ exactly 1 so

We face vanishing or exploding gradient problem.

One method to solve this problem, is the use of neural network architecture that is specifically designed to ensure that derivative of recurrent function is 1. This neural network architecture is LSTM.

The most fundamental idea behind LSTM is that in addition to standard hidden state h_t it also has memory cell c_t for which gradient $\frac{dc_t}{dc_{t-1}}$ is exactly 1. Because the gradient is 1 it do not suffer from vanishing or exploding gradient problem and can capture long term dependency efficiently.