

# Deep Learning Notes

Harshil Pandey

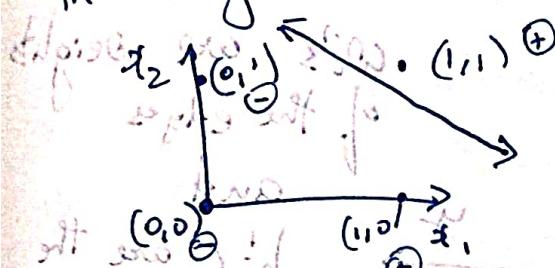
2025

## Contents

1. Basics	2
2. Feed Forward Neural Network	8
3. Momentum Based Gradient Descent and Variants	21
4. Learning Rate Schemes	31
5. Bias and Variance	33
6. Ways to improve the Network Performance	36
7. Activation Functions	41
8. Convolutional Neural Network (CNN)	45
9. Ways to Visualize a CNN	51
10. Joint Distributions and Bayesian Networks	57
11. Undirected Graphical Models	61
12. Latent Variables	63
13. Restricted Boltzman Machines (RBMs)	65
14. Variational Autoencoders (VAEs)	76
15. NADE and MADE	79
16. Generative Adversarial Networks (GANs)	82

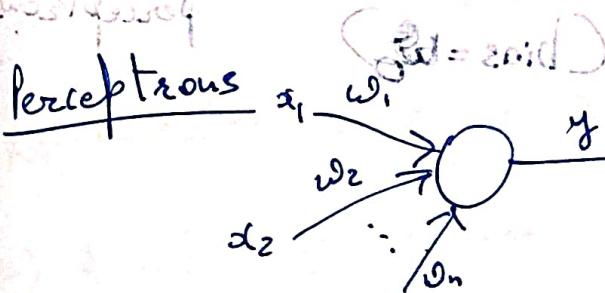
## Linearly Separable Functions

A linearly separable function is a boolean function that can be separated into different classes using a hyperplane in a high-dimensional space.



AND Function is linearly separable.

XOR function is not a linearly separable function.



a)  $w_0$  is called the bias. If  $w_0 > 0$ , we support the class more.

$y = 1$  if  $w_0 + \sum_{i=1}^n w_i x_i > 0$ . We say that the neuron fires if  $w_0 + \sum_{i=1}^n w_i x_i > 0$ .

$\Rightarrow$  Any boolean function can be implemented using a network of perceptrons.

Assume that True = +1 and False = -1. There are 4 possible combinations for 2 inputs.

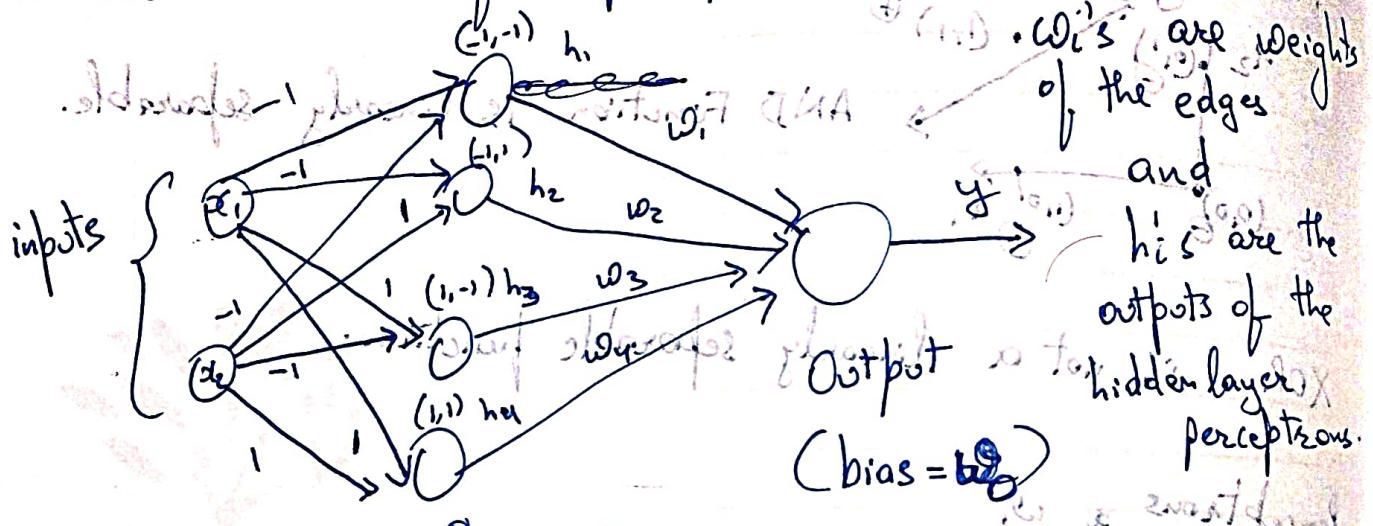
$\Rightarrow$  For 2 inputs, the combinations can be  $2^2 = 4$ . Hence 2 boolean functions can be made for 2 inputs. Each combination of input can take output of -1 or 1. If input can take result of 2 boolean functions: therefore so for n inputs total possible boolean function:  $2^n$ .

$\Rightarrow$  Example for 2 inputs.

total possible input combinations:

$(1, 1), (-1, 1), (1, -1), (-1, -1)$

let us make four perceptrons:



$w_i$ 's are weights of the edges

and

$h_i$ 's are the outputs of the hidden layer

(bias =  $b$ )

It holds if  
to fire,  $\sum w_i h_i + b \geq 0$  ( $b = -2$ )

The claim is, that this network of perceptrons can represent any boolean function of 2 inputs.

This is because if the input is  $(-1, -1)$  only the top neuron will fire, as its output will be  $(-1) \times (-1) + (-1) \times (-1) - 2 = 0 \geq 0$

No other neuron will fire, because we have designed the network in such a way. Hence for each input only the corresponding neuron will fire. So whatever is the desired output for a particular input, that corresponding weight from the corresponding neuron to the output neuron can be adjusted. For example if the desired output from the input  $(-1, -1)$  is 1, then the ~~second~~ only third neuron will

first, hence  $w_3$  can be adjusted w.r.t. to  $b_0$  such that  $\sum_{i=1}^3 w_i b_i + b_0 \geq 0$ . Since no other neuron fires,  $w_3$  can be adjusted independently as it will be the only one contributing to  $\sum_{i=1}^3 w_i b_i + b_0$  when input is  $(1, -1)$ , also  $w_3$  will never contribute whenever input is other than  $(1, -1)$  so it will not matter what  $w_3$  is in those cases.

$\Rightarrow$  When there are 3 inputs, 8 possible input combinations are there. So we use 8 perceptrons in the hidden layer with bias = -3, each firing only to 1 possible input.

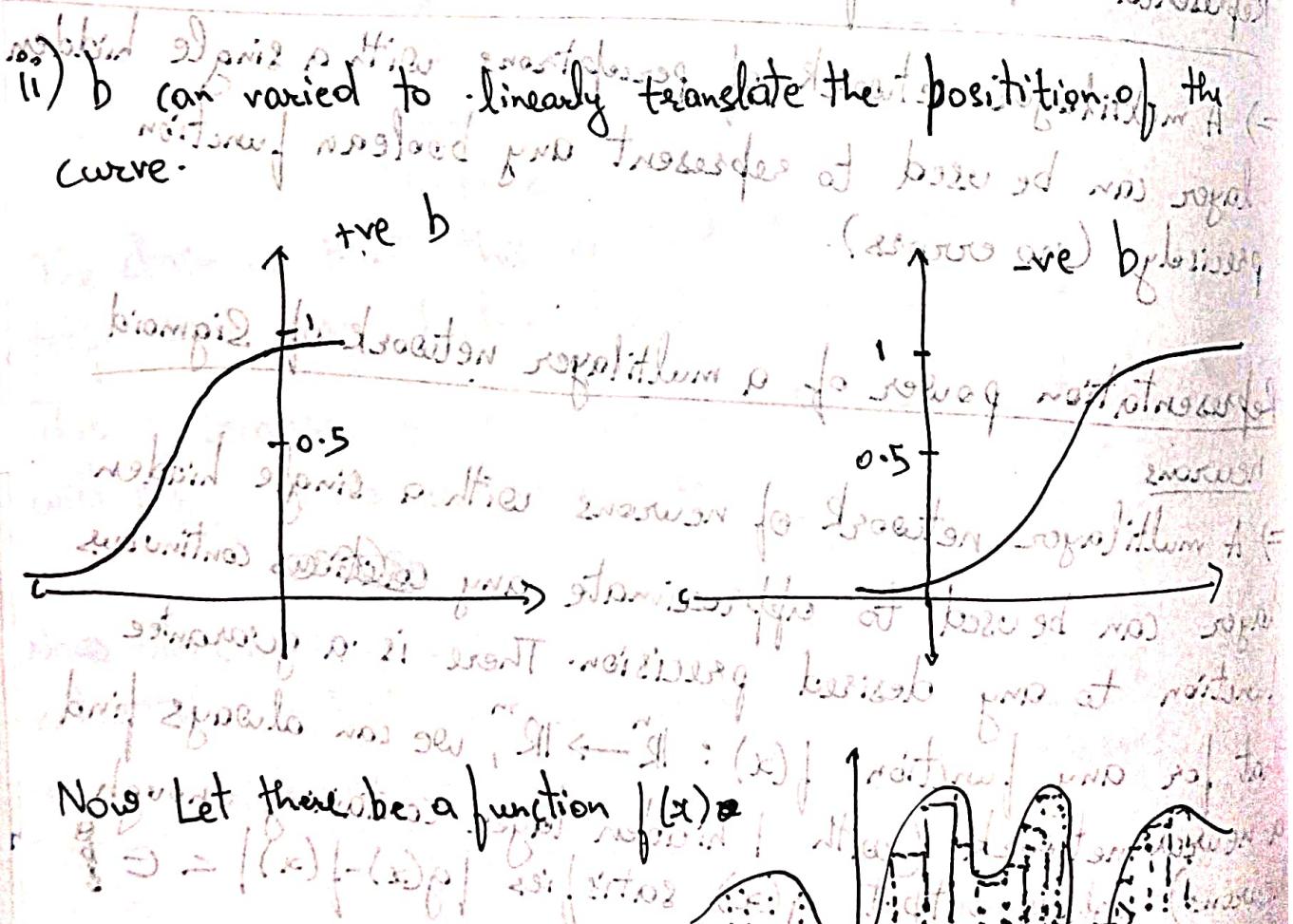
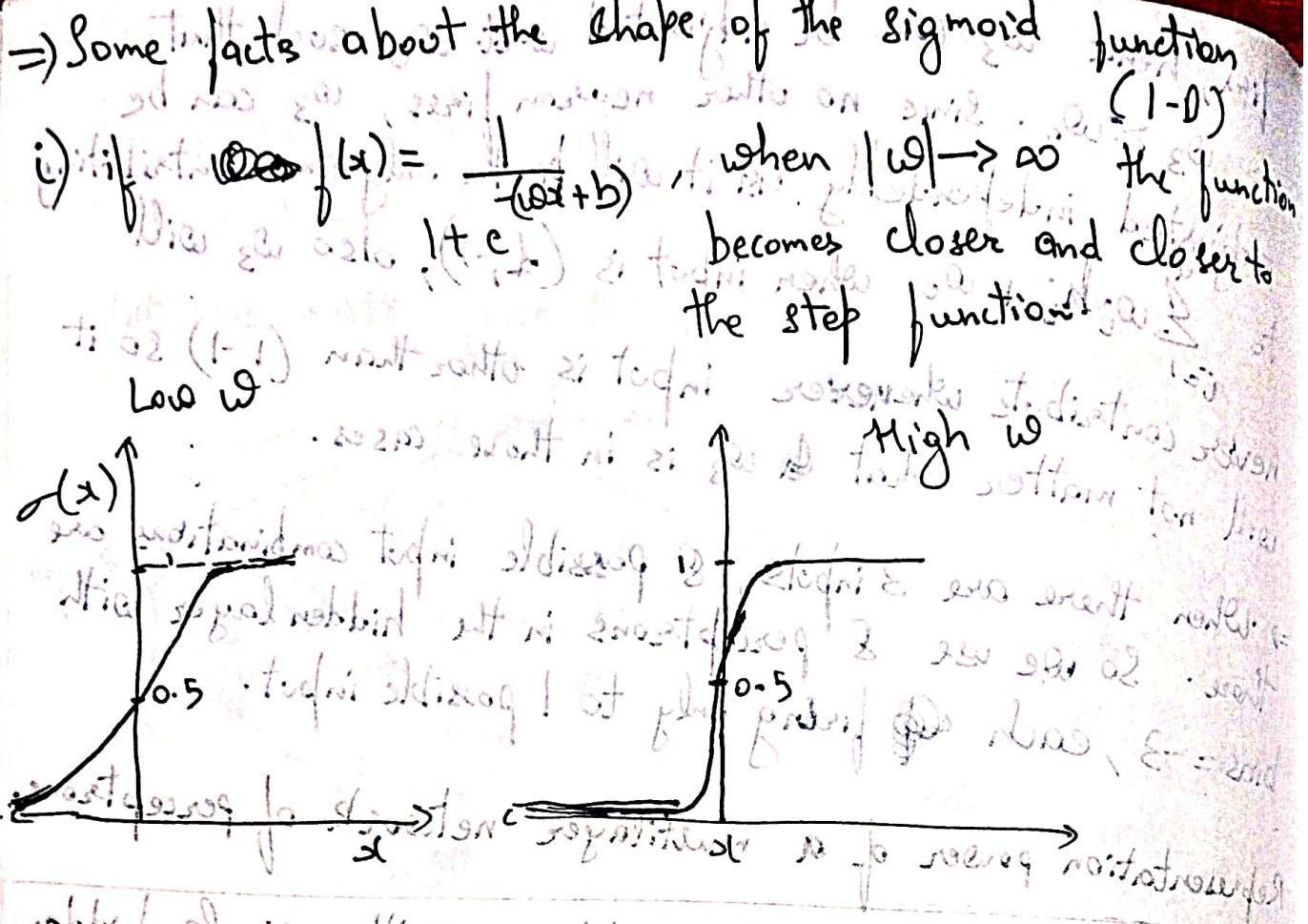
### Representation power of a multilayer network of perceptrons:

$\Rightarrow$  A multilayer network of perceptrons with a single hidden layer can be used to represent any boolean function precisely (no errors).

### Representation power of a multilayer network of Sigmoid neurons

$\Rightarrow$  A multilayer network of neurons with a single hidden layer can be used to approximate any continuous function to any desired precision. There is a guarantee that for any function  $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , we can always find a neural network (with 1 hidden layer containing enough neurons): whose output  $g(x)$  satisfies  $|g(x) - f(x)| < \epsilon$

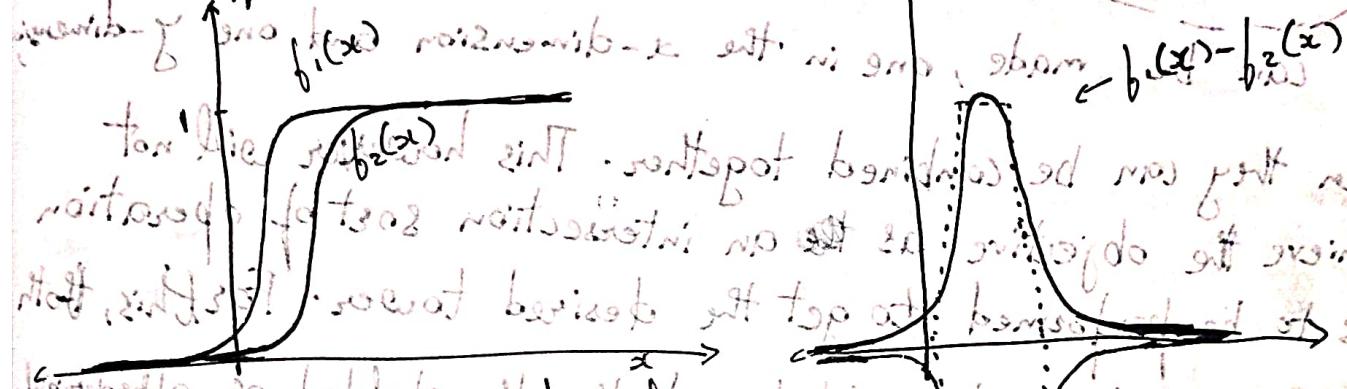
for any  $\epsilon > 0$



This can be approximated through equal width rectangles.

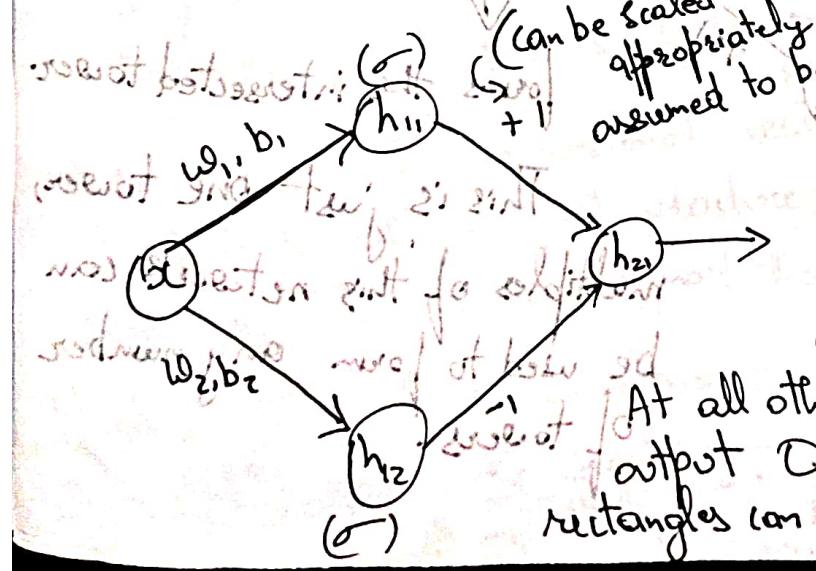
If we introduce more rectangles by reducing the width of the rectangles, then the approximation is more accurate. The question is how to create these rectangle like structures using sigmoid neurons.

Consider two sigmoid functions with very high  $w$  and small difference in bias. Now consider  $f_1(x) - f_2(x)$ .



By adjusting the bias values  $b_1$  and  $b_2$  in  $f_1(x)$  and  $f_2(x)$ , we can manipulate the position of the rectangle. By adjusting the width of the rectangle, we can also manipulate it to translate the like structure linearly.

The height of the rectangle can be manipulated by scaling the output of the sigmoid appropriately. Combining this into a neural network:

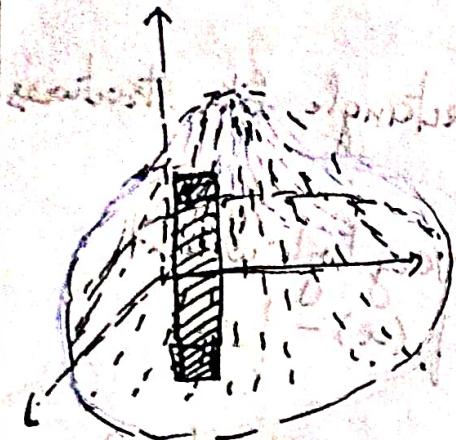


(can be scaled appropriately to achieve the correct height, assumed to be 1 for simplicity)

Notice that by adjusting  $w$ 's and  $b$ 's we can make a very specific rectangle at desired position and height.

At all other inputs, the network will output 0 and won't fire. Any number of rectangles can be made to approximate the function by using the multiples of this network.

For 2-D case each dimension can have its own tower. Suppose the tower has to be made between the area :-



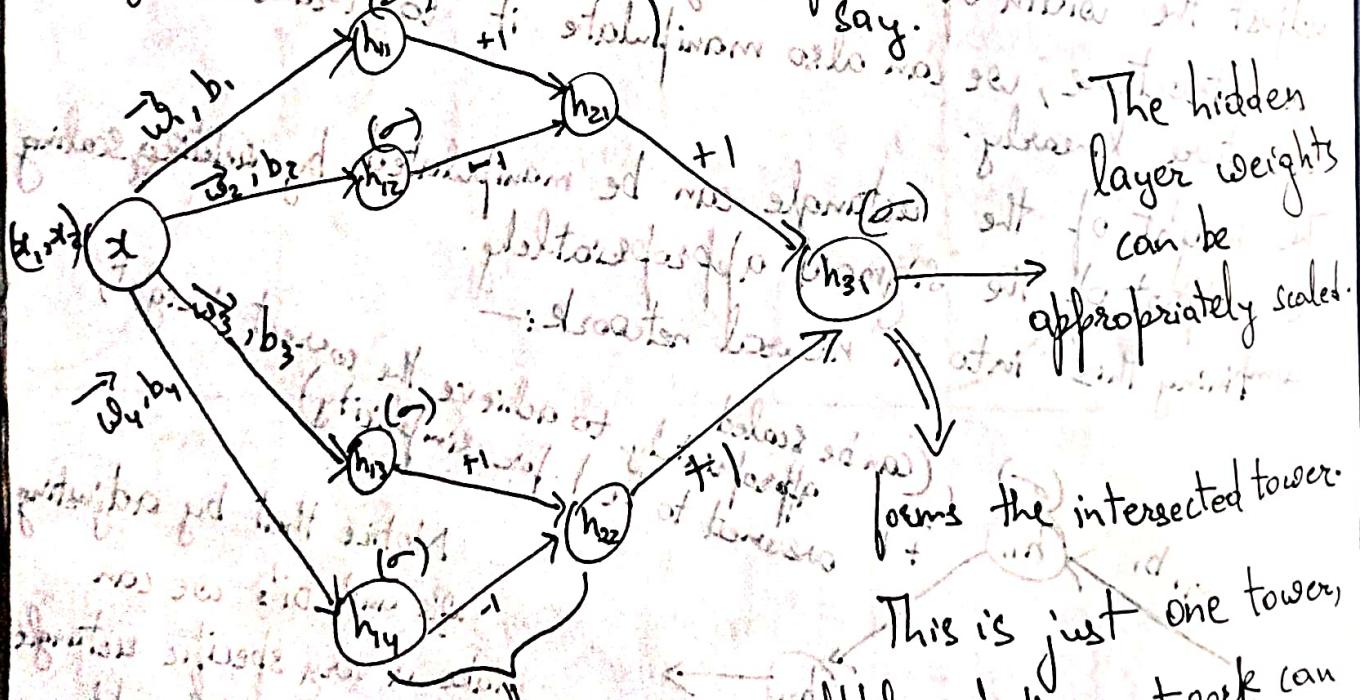
$$\left[ (x_{\min} - x_{\max}) \right] \quad \left[ (y_{\min} - y_{\max}) \right]$$

then two independent towers in the

can be made, one in the  $x$ -dimension and one  $y$ -dimension,

then they can be combined together. This however will not achieve the objective as an intersection sort of operation has to be performed to get the desired tower. For this, both

towers output can be added and then thresholded or appropriately passed through a sigmoid neuron to only fire if the combined value is greater than some value.



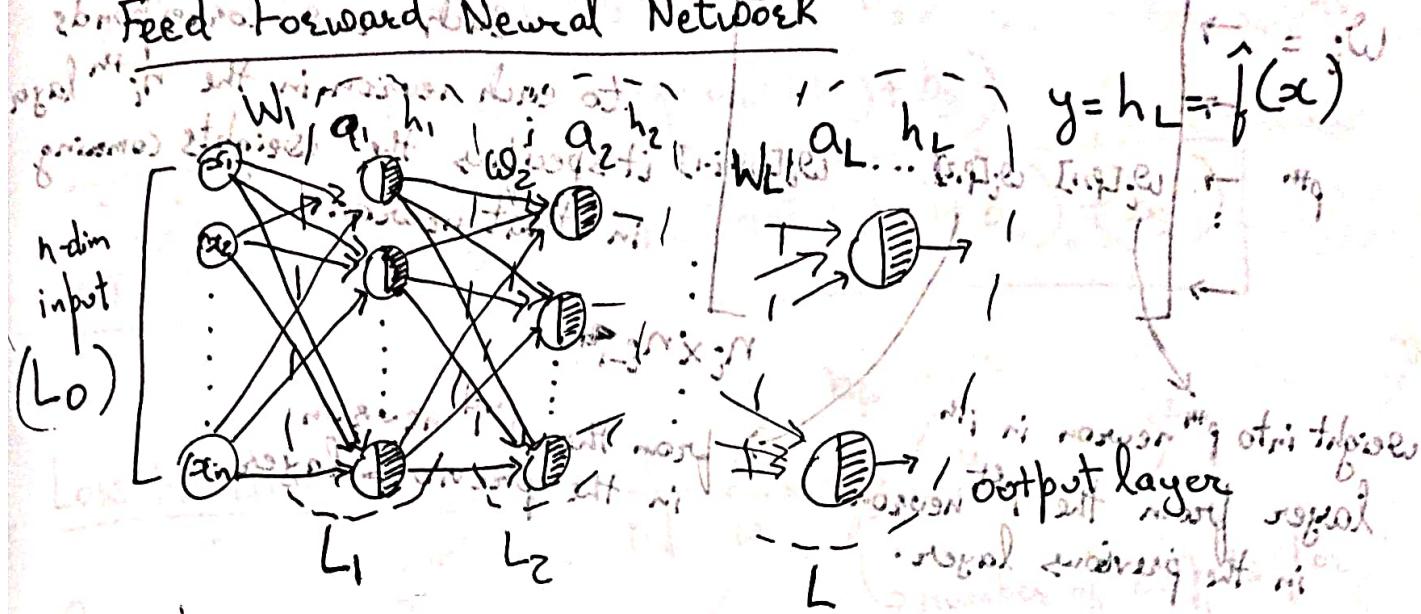
The hidden layer weights can be appropriately scaled.

forms the intersected tower.

This is just one tower, multiples of this network can be used to form any number of towers.

forms the tower in the 2nd dimension say.

## Feed Forward Neural Network



$\Rightarrow$  The input to the network is an  $n$ -dimensional vector  $x \in \mathbb{R}^n$ . The input Layer is called the  $0^{\text{th}}$  Layer.

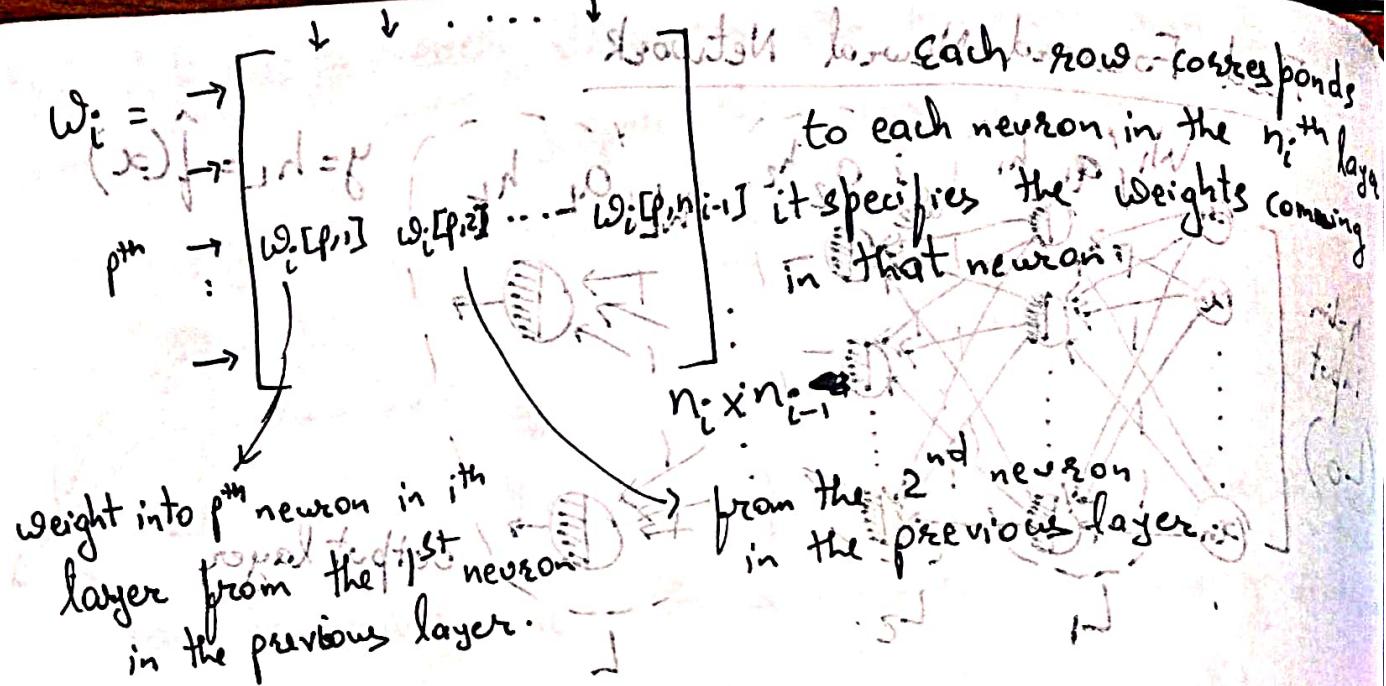
$\Rightarrow$  The network contains  $L-1$  hidden layers, each having  $n_i$  neurons.  $n_0 = n$  (input).  $i$  represents the layers index.

$\Rightarrow$  Each neuron in the hidden layers and output layer can be split into two parts: pre-activation and activation.

$\Rightarrow$   $q_i$  is an  $n$ -dimensional vector such that the  $j^{\text{th}}$  component of  $q_i$  is the input given to the  $j^{\text{th}}$  neuron of the  $i^{\text{th}}$  layer  $L_i$ .

$\Rightarrow$   $h_i$  is an  $n$ -dimensional vector such that the  $j^{\text{th}}$  component of  $h_i$  is the output produced by the  $j^{\text{th}}$  neuron of the  $i^{\text{th}}$  layer  $L_i$ . If neuron is sigmoid then it is equivalent to  $f_j(w_i x + b)$ .

$\Rightarrow$   $w_i \in \mathbb{R}^{n \times n_{i-1}}$  is a matrix corresponding to layer  $i$ .  $w_i[p, q]$  is the weight given to the input from the  $q^{\text{th}}$  neuron in  $L_{i-1}$  layer to  $p^{\text{th}}$  neuron in  $L_i$  layer.



- 3) The pre-activation of the  $i^{th}$  layer will be:
- $$a_i = \sum_{j=1}^{n_{i-1}} w_{i,j} h_{i-1,j} + b_i$$
- ↳ The activation  $a_i$  of the  $i^{th}$  layer is the element wise sum of weighted inputs plus bias.
- ↳ The activation of the  $i^{th}$  layer ( $h_i$ ) is the element wise function of  $a_i$ , for example an element wise sigmoid:-
- $$h_i = \sigma(a_i)$$
- ↳ The activation at the Output Layer is given by:
- $$f(x) = h_L = O(a_L)$$
- ↳ Writing it as a function of  $x$  (assume only 2 hidden layers are there),
- $$f(x) = O(a_3)$$
- $$= O(w_3 h_2 + b_3)$$
- $$= O(w_3 g(a_2) + b_3)$$

$$\hat{f}(x) := O(\omega_3 g(\omega_2 h(x + b_2) + b_3))$$

$$= O(\omega_3 g(\omega_2 g(a) + b_2) + b_3)$$

$$\boxed{\hat{f}(x) = O(\omega_3 g(\omega_2 g(\omega_1 x + b_1) + b_2) + b_3)} \quad (9)$$

(total cost)  $\rightarrow$  works

$$\text{Loss Function: } L(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{n_L} (\hat{y}_{ij} - y_{ij})^2$$

To lost training with N data points.  $\hat{y}_{ij}$ : adding all outputs for each data point.

$\{w_1, w_2, \dots, w_L\}$ : number of outputs for each data point.

$\{b_1, b_2, \dots, b_L\}$ : offset for input data points.

[This is suitable for regression type problems.]

Output Function for regression problems: We would like to allow all real values as outputs so just a linear transformation can work nicely.

$h_L = O(a_L) = w_L a_L + b_L$ , ( $w_L$  and  $b_L$  are not the same as  $w_1$  and  $b_1$ )

(e) For Multiclass Classification: The true labels are  $\{0, 1\}$

where only a single 1 is there in the output vector for 1 datapoint. all other are zero. can be thought of as a PMF with all its mass on the correct class. For modelling such PMF's the most appropriate output activation is the softmax function:

$$\hat{y}_j = O(a_L)_j = \frac{e^{a_{Lj}}}{\sum_{k=1}^n e^{a_{Lk}}} = \left\{ \begin{array}{l} \text{the } j^{\text{th}} \text{ component of an} \\ \text{output vector for 1 data} \\ \text{point. } \hat{y}_j \in (0, 1) \end{array} \right.$$

In this case, the more appropriate loss function is the cross-entropy loss  $L(\theta) = -\sum_{j=1}^{n_r} y_j \log(\hat{y}_j)$  (Loss for 1 data point).

$$\Rightarrow L(\theta) = -\log(\hat{y}_t) \quad \text{where } y_t = 1 \quad (\text{true label})$$

$\hat{y}_j$  represents the probability of the data point that it belongs to the  $j^{\text{th}}$  class.

Gradient Computations (Cross entropy loss and softmax output neurons)

Let for some data point, the correct class is  $t^{\text{th}}$  class. Then the output vector for that point will be  $e_t = [0 \dots 0, 1, 0]$

The loss corresponding to this point is:

$$L(\theta) = -\log(\hat{y}_t)$$

Let's calculate  $\nabla_{\theta} L(\theta)$ , gradient with respect to the output layer

$$\nabla_{\theta} L(\theta) = \nabla_{\hat{y}} L(\theta)$$

Let us start with  $\nabla_{\hat{y}} L(\theta)$  - the  $t^{\text{th}}$  element of  $\nabla_{\hat{y}} L(\theta)$ .

$$\Rightarrow \frac{\partial L(\theta)}{\partial \hat{y}_j} = \frac{\partial}{\partial \hat{y}_j} (-\log(\hat{y}_t)) = 1/t \times -\frac{1}{\hat{y}_t}$$

$$\left( \frac{\partial L(\theta)}{\partial \hat{y}_j} \right) = \begin{cases} \frac{1}{\hat{y}_t} & \text{if } j=t \\ 0 & \text{if } j \neq t \end{cases}$$

$$\nabla_{\hat{y}} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial \hat{y}_1} \\ \frac{\partial L(\theta)}{\partial \hat{y}_2} \\ \vdots \\ \frac{\partial L(\theta)}{\partial \hat{y}_{n_L}} \end{bmatrix} = -\frac{1}{\hat{y}_t} \begin{bmatrix} \frac{1}{1-(t-t)} \\ \frac{1}{1-(2-t)} \\ \vdots \\ \frac{1}{1-(n_L-t)} \end{bmatrix}$$

$$\frac{\partial L(\theta)}{\partial \hat{y}_t} = -\frac{1}{\hat{y}_t} \times e_t$$

$$\boxed{\nabla_{\hat{y}} L(\theta) = -\frac{1}{\hat{y}_t} \times e_t}$$

Now let's calculate the gradient based on the pre-activation layer of the output.  $\nabla_{a_L} L(\theta)$ .

For that we know that:-  $\hat{y}_j = \frac{e^{a_{Lj}}}{\sum_{k=1}^n e^{a_{Lk}}} \quad \{ \text{so } t \text{ max.} \}$

$$L(\theta) = -\log(\hat{y}_t)$$

$$\frac{\partial L(\theta)}{\partial a_{Lj}} = \cancel{\frac{\partial (-\log(\hat{y}_t))}{\partial a_{Lj}}} = \frac{\partial (-\log(\hat{y}_t))}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial a_{Lj}}$$

$\Downarrow$

$$\frac{\partial L(\theta)}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial a_{Lj}}$$

$$= -\frac{1}{\hat{y}_t} \cdot \frac{\partial}{\partial \hat{y}_t} \left( \frac{e^{a_{L,t}}}{\sum_{k=1}^n e^{a_{Lk}}} \right)$$

$$= -\frac{1}{\hat{y}_t} \cdot \frac{\partial}{\partial \hat{y}_t} \left( \frac{\sum_{k=1}^n e^{a_{Lk}} \cdot \frac{\partial (a_{L,t})}{\partial a_{Lj}}}{\left( \sum_{k=1}^n e^{a_{Lk}} \right)^2} \right)$$

$$= -\frac{1}{y_t} \times \left( \frac{\sum_{k=1}^{n_L} e^{\alpha_{Lk}} \alpha_{Lt} \mathbb{1}(t=j) \frac{(0.16)}{0.06} e^{\alpha_{Lj}}}{\left( \sum_{k=1}^{n_L} e^{\alpha_{Lk}} \right)^2} \right)$$

$$= -\frac{1}{y_t} \times \left( \frac{\mathbb{1}(t=j) e^{\alpha_{Lj}}}{\sum_{k=1}^{n_L} e^{\alpha_{Lk}}} - \frac{e^{\alpha_{Lj}}}{\sum_{k=1}^{n_L} e^{\alpha_{Lk}} \times \left( \sum_{k=1}^{n_L} e^{\alpha_{Lk}} \right)} \right)$$

Belongs to the following

$$= -\frac{1}{y_t} \left( \frac{\mathbb{1}(t=j) \hat{y}_j}{\sum_{k=1}^{n_L} e^{\alpha_{Lk}}} \right)$$

This belongs to the following equation  
with respect to the first term of the objective function.

$$= -\left( \mathbb{1}(t=j) \hat{y}_j \right)$$

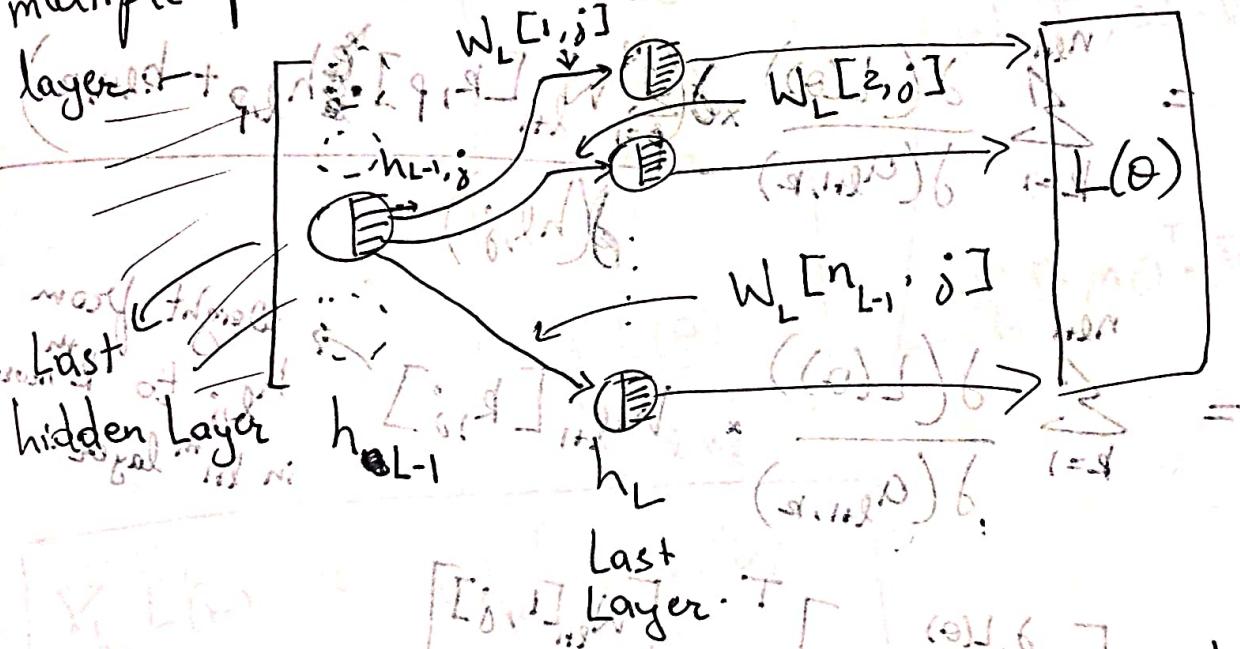
$$\Rightarrow \frac{\partial L(\theta)}{\partial \alpha_{Lj}} = -\left( \mathbb{1}(t=j) - \hat{y}_j \right) \frac{(0.16)}{0.06} = \frac{(0.16)}{0.06}$$

Therefore:  $\nabla_{\alpha_L} L(\theta) = - \begin{bmatrix} \mathbb{1}(t=1) - \hat{y}_1 \\ \mathbb{1}(t=2) - \hat{y}_2 \\ \vdots \\ \mathbb{1}(t=n_L) - \hat{y}_{n_L} \end{bmatrix}$

$$\boxed{\nabla_{\alpha_L} L(\theta) = -(e_t - \hat{y})}$$

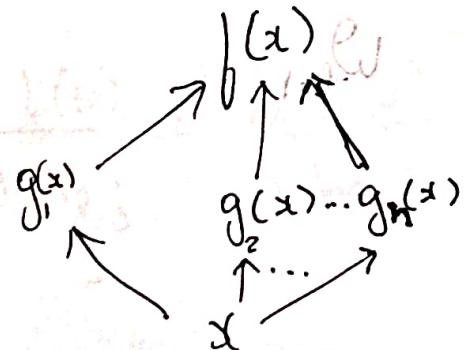
We now have the formulas for gradients of the ~~next~~ output layer. This is straightforward as the output layer directly goes into the loss function.

This is not true however, for any layer ~~except~~ as the output of hidden layer affects the loss function through multiple paths. Let us take a neuron in the last hidden layer.



We can see that output of  $h_{L-1, j}$  <sup>th</sup> neuron of the last hidden layer affects the path through  $L$  paths. Hence we use chain rule for multiple path to calculate

$\frac{\partial L(\theta)}{\partial h_{L-1, j}}$  <sup>th</sup> neuron of the last hidden layer.   
 l is for the layer number and j is for the neuron number.



Let  $f(x)$  be a function of  $x$  such that it can be separated as  $g_1(x), g_2(x)$  and  $\dots g_n(x)$ . Then

$$\frac{\partial f}{\partial x} = \sum_{i=1}^n \frac{\partial f}{\partial g_i(x)} \times \frac{\partial g_i(x)}{\partial x}$$

In our case,  $f(x) = L(\theta)$ ,  $g_i(x) = \alpha_{l+1, k}$   
 all we have to find is  $\frac{\partial L(\theta)}{\partial h_{l,j}}$   $\rightarrow$    
 we know  $x = h_{l,j}$  is <sup>in</sup> <sup>netto</sup> <sup>pre-activation</sup> of the next layer.  
 $\frac{\partial L(\theta)}{\partial h_{l,j}} = \sum_{k=1}^{n_{l+1}} \frac{\partial L(\theta)}{\partial \alpha_{l+1, k}} \frac{\partial \alpha_{l+1, k}}{\partial h_{l,j}}$   $\rightarrow$   $\frac{\partial L(\theta)}{\partial h_{l,j}} = \sum_{k=1}^{n_{l+1}} \alpha_{l+1, k} \frac{\partial \alpha_{l+1, k}}{\partial h_{l,j}}$

$$\frac{\partial L(\theta)}{\partial h_{l,j}} = \sum_{k=1}^{n_l} \frac{\partial L(\theta)}{\partial a_{l+1,k}} \frac{\partial a_{l+1,k}}{\partial h_{l,j}}$$

$$= \sum_{k=1}^{n_{l+1}} \frac{\partial L(\theta)}{\partial \alpha_{l+1,k}} \times \left( \sum_{p=1}^{n_l} w_{l+1}^{[k,p]} h_{l,p} + b_{l+1,k} \right)$$

$$= \sum_{k=1}^{n_{l+1}} \frac{\delta(L(\theta))}{\delta(a_{l+1,k})} \times W_{l+1}[k, i]$$

Weight from  
h<sub>l+1</sub><sup>i</sup> to h<sub>k</sub><sup>i</sup>  
in l+1<sup>th</sup> layer

•  $\text{grad} = \left[ \frac{\partial L(\theta)}{\partial \theta_1}, \dots, \frac{\partial L(\theta)}{\partial \theta_n} \right]$   $\rightarrow$   $W_{l+1}^{[1], \theta}$   
 - 2. If no mistakes in  $\theta$ :  
 -  $\frac{\partial L(\theta)}{\partial \theta_l}$   $\rightarrow$   $W_{l+1}^{[2], \theta}$   $\rightarrow$   $W_{l+1}^{[n], \theta}$   $\rightarrow$  final output  
 - mistakes at step  $\theta_l$ :  $\rightarrow$   $W_{l+1}^{[1], \theta}$   $\rightarrow$   $W_{l+1}^{[2], \theta}$   $\rightarrow$   $W_{l+1}^{[n], \theta}$   $\rightarrow$  see mistake

$\frac{\partial L(\theta)}{\partial \theta_j}$  is equal to sum of all the  $j^{th}$  column of  $A^T A$  multiplied by  $\theta_j$ .

$$\nabla L(\theta) = \frac{\partial}{\partial \theta} \sum_{i=1}^m f_i(\theta)$$

$$\Rightarrow \frac{\partial L(\theta)}{\partial h_{l,j}} = W_{l+1}^T [ \cdot, j ]^T \cdot \nabla_{a_{l+1}} L(\theta) = (\theta)_j \nabla_{a_{l+1}} L(\theta)$$

(s.t.)  $\theta = [ \cdot ]_6$

Here

$$\nabla_{h_l} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial h_{l,1}} \\ \frac{\partial L(\theta)}{\partial h_{l,2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{l,n_l}} \end{bmatrix}$$

$$= W_{l+1}^T [ \cdot, 1 ]^T \nabla_{a_{l+1}} L(\theta)$$

$$= W_{l+1}^T [ \cdot, 2 ]^T \nabla_{a_{l+1}} L(\theta)$$

$$\vdots$$

$$= W_{l+1}^T [ \cdot, n_l ]^T \nabla_{a_{l+1}} L(\theta)$$

What happens if there are multiple layers?  $L(\theta) = \dots \circ (\theta)_1 \circ \dots \circ (\theta)_L$

also  $(\theta)_i \circ \theta = (\theta)_i$

$$\nabla_{h_l} L(\theta) = W_{l+1}^T \nabla_{a_{l+1}} L(\theta)$$

Now let us look at the  $\nabla_{a_l} L(\theta)$ . We know that

~~$h_l = g(a_l)$~~

$h_{l,j} = g(a_{l,j})$ 

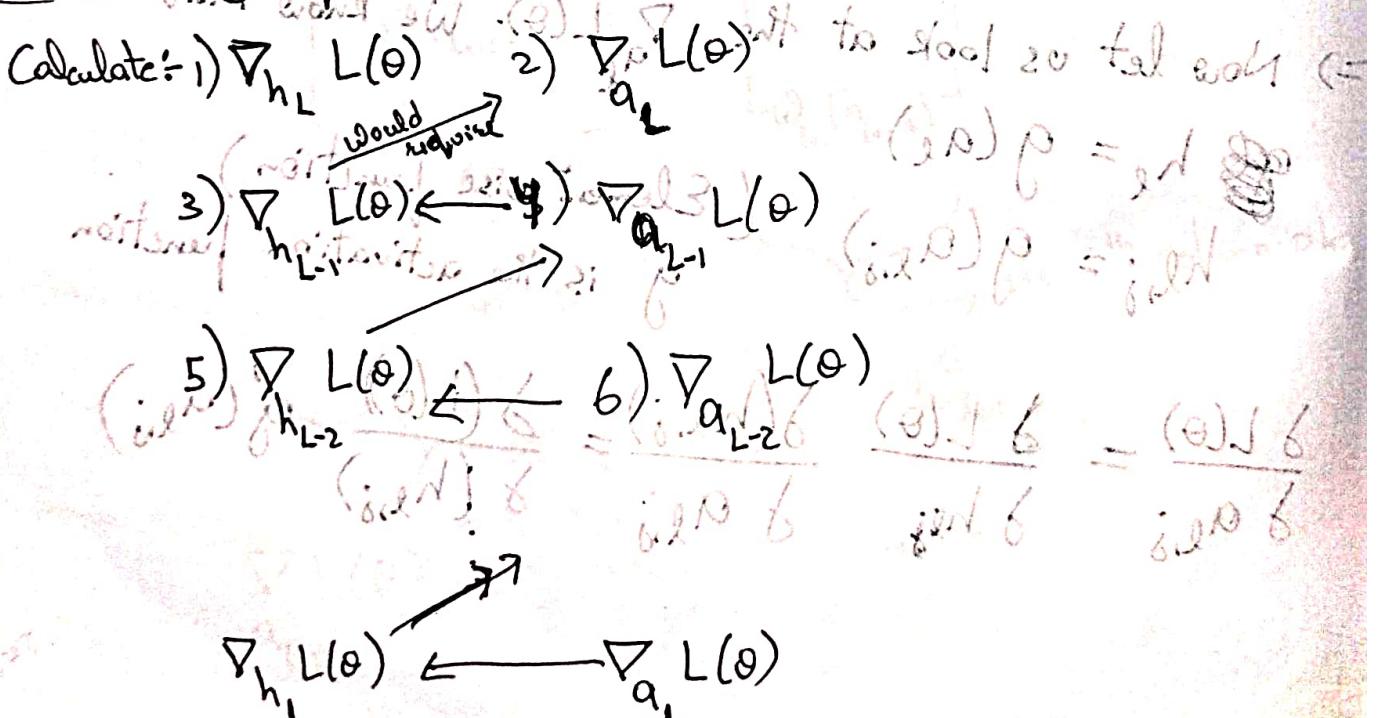
$g$  is the activation function

$$\frac{\partial L(\theta)}{\partial a_{l,j}} = \frac{\partial L(\theta)}{\partial h_{l,j}} \quad \frac{\partial (h_{l,j})}{\partial a_{l,j}} = \frac{\partial (L(\theta))}{\partial (h_{l,j})} g'(a_{l,j})$$

$$\nabla_{a_e} L(\theta) = \left[ \begin{array}{c} \frac{\partial L(\theta)}{\partial a_{e,1}} \\ \frac{\partial L(\theta)}{\partial a_{e,2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{e,n_e}} \end{array} \right] = \left[ \begin{array}{c} \frac{\partial L(\theta)}{\partial h_{e,1}} \cdot g'(a_{e,1}) \\ \frac{\partial L(\theta)}{\partial h_{e,2}} \cdot g'(a_{e,2}) \\ \vdots \\ \frac{\partial L(\theta)}{\partial h_{e,n_e}} \cdot g'(a_{e,n_e}) \end{array} \right]$$

$\nabla_{a_e} L(\theta) = \nabla_{h_e} L(\theta) \odot [\dots, g'(a_{e,i}) \dots]$  Element wise multiplication of  $\nabla_{h_e} L(\theta)$  and  $g'(a_e)$  also called as Hadamard Product

### Chronology of Gradient Computation



$\Rightarrow$  Computing gradients for the parameters  $\nabla_{W_k} L(\theta)$   $n_k \times n_{k-1}$   
 We know that :-

$$a_l = W_l h_{l-1} + b_l$$

$$\text{and } \nabla_b L(\theta)$$

$$\frac{\partial L(\theta)}{\partial W_l[j, k]}$$

$$\frac{\partial L(\theta)}{\partial W_l[i, k]}$$

Note that  $a_{l,j}$  only comes in the computation of  $a_{l,j}$ . Hence

$W_l[j, k]$  affects  $L(\theta)$  only

through 1 path.

Hence:-

$$\frac{\partial L(\theta)}{\partial W_l[j, k]} = \frac{\partial L(\theta)}{\partial a_{l,j}} \times \frac{\partial a_{l,j}}{\partial W_l[j, k]}$$

Weight to  $j^{th}$  neuron from  $k^{th}$  neuron

$$= \frac{\partial L(\theta)}{\partial a_{l,j}} \times \frac{\partial}{\partial W_l[j, k]} \left( \sum_{p=1}^{n_{l-1}} W_l[j, p] \cdot h_{l-1, p} + b_{l,j} \right)$$

$$\frac{\partial L(\theta)}{\partial W_l[j, k]} = \frac{\partial L(\theta)}{\partial a_{l,j}} \times h_{l-1, k}$$

$$\Rightarrow \nabla_{W_l} L(\theta) =$$

$$\nabla_{W_l} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial W_l[1, 1]} & \frac{\partial L(\theta)}{\partial W_l[1, 2]} & \dots & \frac{\partial L(\theta)}{\partial W_l[1, n_{l-1}]} \\ \frac{\partial L(\theta)}{\partial W_l[2, 1]} & \frac{\partial L(\theta)}{\partial W_l[2, 2]} & \dots & \frac{\partial L(\theta)}{\partial W_l[2, n_{l-1}]} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L(\theta)}{\partial W_l[n_{l-1}, 1]} & \frac{\partial L(\theta)}{\partial W_l[n_{l-1}, 2]} & \dots & \frac{\partial L(\theta)}{\partial W_l[n_{l-1}, n_{l-1}]} \end{bmatrix}$$

$$\nabla_{W_l} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{l,1}} \cdot h_{l-1,1} & \frac{\partial L(\theta)}{\partial a_{l,2}} \cdot h_{l-1,2} & \dots & \frac{\partial L(\theta)}{\partial a_{l,n_{l-1}}} \cdot h_{l-1,n_{l-1}} \\ \frac{\partial L(\theta)}{\partial a_{l,1}} \cdot h_{l-1,1} & \frac{\partial L(\theta)}{\partial a_{l,2}} \cdot h_{l-1,2} & \dots & \frac{\partial L(\theta)}{\partial a_{l,n_{l-1}}} \cdot h_{l-1,n_{l-1}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L(\theta)}{\partial a_{l,1}} \cdot h_{l-1,1} & \frac{\partial L(\theta)}{\partial a_{l,2}} \cdot h_{l-1,2} & \dots & \frac{\partial L(\theta)}{\partial a_{l,n_{l-1}}} \cdot h_{l-1,n_{l-1}} \end{bmatrix}$$

$$\Rightarrow \nabla_{W_L} L(\theta) = \begin{bmatrix} \frac{\partial L(\theta)}{\partial a_{L,1}} \\ \frac{\partial L(\theta)}{\partial a_{L,2}} \\ \vdots \\ \frac{\partial L(\theta)}{\partial a_{L,N_L}} \end{bmatrix} = \begin{bmatrix} h_{L-1,1} & h_{L-1,2} & \dots & h_{L-1,N_{L-1}} \end{bmatrix}^T$$

(a) J 6

$$\boxed{\nabla_{W_L} L(\theta) = \nabla_{a_L^T} L(\theta) \cdot h_{L-1}^T}$$

(a) J 6

Simpler argument  
as  $\frac{\partial L(\theta)}{\partial a_{L,i}}$   
is just  $\frac{\partial L(\theta)}{\partial b_L}$

For  $\nabla_{a_L^T} L(\theta)$ , notice that  $\frac{\partial L(\theta)}{\partial a_{L,i}} = \frac{\partial L(\theta)}{\partial a_{L,0}} \times \frac{\partial a_{L,0}}{\partial a_{L,i}}$

$$\text{So } \boxed{\nabla_{a_L^T} L(\theta) = \nabla_{a_L^T} L(\theta) \cdot \frac{\partial a_{L,0}}{\partial a_{L,i}}}$$

(a) J 6

$$\Rightarrow \sigma(x) = \frac{1}{1+e^{-x}} \quad \sigma'(x) = \sigma(x)(1-\sigma(x))$$

(a) J 6

$$\Rightarrow \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad \tanh'(x) = 1 - (\tanh(x))^2$$

(a) J 6

$$\frac{e^x - e^{-x}}{e^x + e^{-x}} \cdot \frac{e^x + e^{-x}}{e^x + e^{-x}}$$

(a) J 6

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

(a) J 6

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

(a) J 6

## Contours

Suppose there is a function from  $\mathbb{R}^2 \rightarrow \mathbb{R}$ ,  $z = f(x, y)$ , then the contour map may look like:-

Each contour ring states the  $(x, y)$  values where the  $f(x, y)$  is same.

$\Rightarrow$  Wherever the contour rings are well located, the function has a steep curvature.  
 $\Rightarrow$  Wherever the contour rings are further apart, the function has a gentle slope/curvature.

The logic is that, in areas of steep curvature, a small change in  $x$  and  $y$  values may bring a larger change in  $f(x, y)$ . Similarly in areas of gentle slopes, a large change is required to bring the same amount of change in  $f(x, y)$ .

$\Delta z$  change is achieved in small changes in  $x$  and  $y$ , steep curvature, contour rings are closer.

$\Delta z$  change is achieved in larger changes in  $x$  and  $y$ , gentle slopes, contour rings are further apart.

$\Rightarrow$  When gradient descent is applied, it makes tiny progress in areas where contour rings are far apart as the slope is small. It makes larger steps in areas where contour rings are closer as the slope is larger.

Rapid progress is made in areas where contour rings are closer.

## Momentum Based Gradient Descent

The small progress made in the flat areas causes the gradient descent to converge slowly. But if we are repeatedly asked to move in the same/similar directions, it makes sense to take bigger steps in that direction. So instead of just guiding our search based on just the gradient, we introduce some historical logic into the search. So the update rule hence becomes:-

where,  $\nabla w_t = \nabla w_t + \beta U_{t-1}$  where  $U_t$  is the combination of the current gradient and historical directions.

where,  
 $U_t = \nabla w_t + \beta U_{t-1}$  we start with  $w_0$  and  $0 \leq \beta < 1$ .  
 $(\beta)$  is called  $\beta$  is called momentum factor.  
 $w_0 = \nabla w_0$ , then  $w_1 = w_0 - \eta \nabla w_0$ .  
 $w_1 = \nabla w_1 + \beta U_0 = \nabla w_1 + \beta \nabla w_0$ .

$$\begin{aligned} \text{So } w_2 &= w_1 - \eta (\nabla w_1 + \beta (\nabla w_0 + \beta \nabla w_0)) \\ &= w_1 - \eta (\nabla w_1 + \beta \nabla w_1 + \beta^2 \nabla w_0) \\ &= w_1 - \eta \sum_{i=0}^{2-1} \beta^i \nabla w_i \quad \text{which is the current gradient} \\ &\quad \text{if we then } \nabla w_1, \text{ and then the least} \\ &\quad \text{importance to } \nabla w_0. \quad \text{so if we make small steps} \\ \therefore w_{t+1} &= w_t - \eta \sum_{i=0}^t \beta \nabla w_i \quad \text{in some direction, after some steps,} \\ &\quad \text{the "momentum" } (U_t) \text{ will grow} \\ &\quad \text{and we will start making bigger} \\ &\quad \text{steps.} \end{aligned}$$

The problem with this strategy is that sometimes the momentum is too large and we overshoot & unnecessarily in some directions. But it certainly converges faster than normal in most situations.

G.D. in most situations is better than momentum based To deal with the overshooting problem of Momentum based gradient descent we look at the Nesterov Accelerated Gradient Descent. Instead of directly combining the historical direction  $U_{t-1}$  with the current gradient to calculate  $U_t$  we check what would happen if  $\nabla w_t$  was not there. Then  $U_t$  will be  $\beta U_{t-1}$  :-  $U_t = \beta U_{t-1}$

$$\text{So } W'_{t+1} = W_t - \eta U_t = W_t - \eta \beta U_{t-1}$$

If we ignore the step size  $\eta$ , then  $W'_{t+1}$  will be  $W_t - \beta U_{t-1}$ . Let us calculate the gradient w.r.t this  $W'_{t+1}$ .

Now let's look at the 2 cases:-

1) Not overshooting:- if we don't overshoot, then  $\nabla W'_{t+1}$  will be in line with the previous directions, so if  $U_t$  is calculated like this:-

$$U_t = \nabla(W'_{t+1}) + \beta U_{t-1} \quad (\text{it would add on the momentum})$$

Since both  $\nabla W'_{t+1}$  and  $\beta U_{t-1}$  are similar directions.

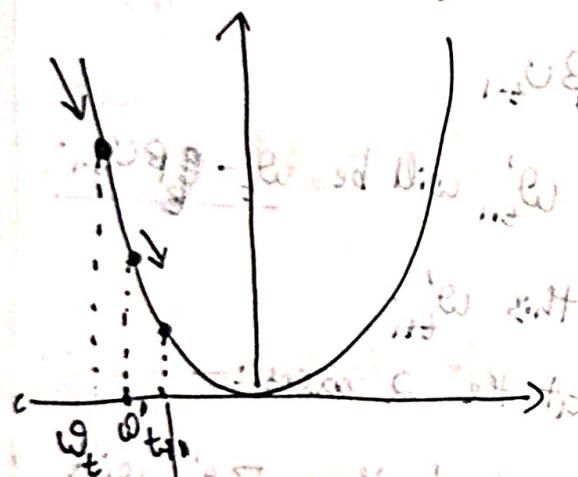
2) Overshoot:- if we overshoot, then  $\nabla W'_{t+1}$  will be in opposite like direction with the previous directions. So if  $U_t$  is calculated like this:-

$U_t = \nabla W_{t+1} + \beta U_{t-1}$ , it would reduce the momentum effect as  $\nabla W_{t+1}$  and  $\beta U_{t-1}$  affect in opposite like directions. Hence this will prevent/reduce overshooting as it would reduce the momentum effect.

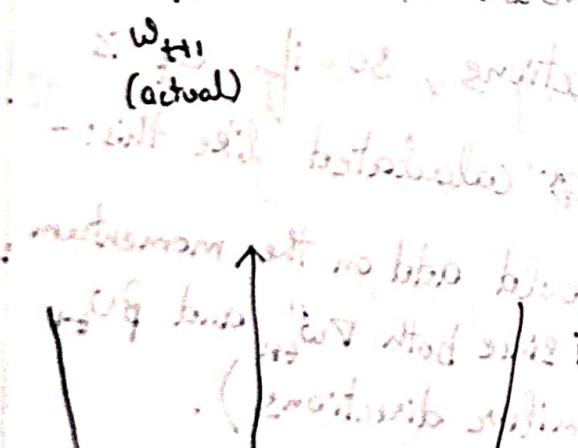
In both cases:-

$$U_t = \nabla(W_t - \beta U_{t-1}) + \beta U_{t-1}$$

$$W_{t+1} = W_t - \eta U_t$$

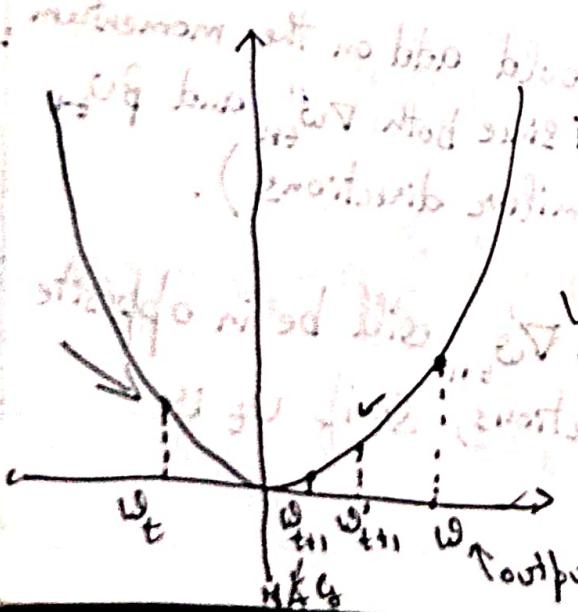


$W_{t+1} = W_t - \beta U_{t-1}$  looking ahead to  $W_{t+1}$ . In this case, looking ahead to  $W_{t+1}$  suggests that we should go along with the momentum and leads us to the next equilibrium.



$W_{t+1} = W_t - \beta U_{t-1}$  looking ahead to  $W_{t+1}$  suggests that we will overshoot and it cancels some momentum to avoid overshoot and leads us to  $W_{t+1}$ .

$W$  would have been the output for Momentum-Based GD. (Overshooting is high)



Output for Momentum Based GD:

so NAG oscillates, lesser than SGD. To stabilize it we can do so:  
1) halve learning rate after every 5 epochs.

One more way to speed up convergence is to set the correct learning rate  $\eta$ . Tips for annealing the learning rate:

- 1) Halve the learning rate after every 5 epochs.
- 2) Halve the learning rate after an epoch if the validation error is more than what it was at the end of the previous epoch. So if the validation error is more at the end of current epoch, then reinitialize the weights to what was learned after previous epoch and run the current epoch again with  $1/2$  the learning rate.
- 3) Exponential Decay: -  $\eta = \eta_0 e^{-kt}$  where  $\eta_0$  and  $k$  are initial hyperparameters and  $t$  is the step number. If  $\eta_0$  or  $k \uparrow$  then the decay is more rapid.

$$4) \frac{1}{t} \text{ decay, } \eta = \frac{\eta_0}{1+kt}$$

### 5) For Momentum algorithms

$$\beta_t = \min\left(1 - 2^{-1 - \log_2(\lfloor \frac{t}{350} \rfloor + 1)}, \beta_{\max}\right), \text{ where}$$

$\beta_{\max}$  is chosen from  $\{0.999, 0.995, 0.99, 0.9, 0\}$ .

- 6) Line Search: try many values of  $\eta$  for each step, choose the one which causes the least loss.

All the methods discussed above change the learning rate for the whole gradient. This means irrespective of the features, the learning rate remains the same for all. But suppose we take the gradient

for the whole dataset and some of the features were important but very sparse. In this case the gradient update will be  $\Delta w_t = w_t - \eta \nabla w_t$ . The benefit of this is that gradients of most of the features will be zero in the gradient.

Naturally for sparse features, the corresponding element will be very small. Hence the updates made for these sparse features will be small as compared to other dense features. At the end of the algorithm (exhausting the iteration limit) the  $w^*$  for the sparse features may not be the optimal value as the updates were small in the iterations. Can we adapt the learning rate for each feature based on their frequency? (sparse features having a relatively higher learning rate).

AdaGrad :- Decay the learning rate for each parameter in proportion to their update history (more updates means more decay).

(larger magnitude)

Update rule for AdaGrad :-  $v_t = v_{t-1} + (\nabla w_t)^2$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \nabla w_t$$

Scalor, as it refers to gradient w.r.t each parameter

1) For dense features,  $v_t$  grows much faster as  $\nabla w$  is large, so the decay in the learning rate is much faster.

2) For sparse features,  $v_t$  grows slowly as  $\nabla w$  is small, so the decay in the learning rate is slow.

$\Rightarrow$  Adagrad slows down near the minimum due to decaying learning rate. This is because after some iterations,  $v_t$  is large and decay is too much, especially for dense features. Moreover, near the minima, the gradients are also small, which when combined with large accumulated  $v_t$  slows down the learning algorithm.  $v_t$  keeps on increasing.

$\Rightarrow$  To avoid this problem, we can use an algorithm that decays the denominator  $/v_t$  in such a way that in the steep areas,  $v_t$  is large but in gentle areas  $v_t$  decreases. (effective learning is small)

This is done in RMS Prop-algorithm.

$$v_t = \beta v_{t-1} + (1-\beta)(\nabla w_t) \quad 0 \leq \beta \leq 1, \text{ typically } 0.9, 0.99$$

let us expand  $v_t$

$$\begin{aligned} v_3 &= 0.9v_2 + 0.1\nabla w_3^2 \\ v_3 &= 0.9(0.9v_1 + 0.1\nabla w_2^2) + 0.1\nabla w_3^2 \\ v_1 &= 0.9v_0 + 0.1\nabla w_1^2 \\ v_0 &= 0.9(0.9v_0 + 0.1\nabla w_0^2) + 0.1\nabla w_1^2 \\ \Rightarrow v_3 &= 0.9(0.9(0.9v_0 + 0.1\nabla w_0^2) + 0.1\nabla w_1^2) + 0.1\nabla w_3^2 \end{aligned}$$

$$v_3 = 0.9 \cdot (0.9 \cdot (0.9 \cdot v_0 + 0.1 \nabla w_0^2) + 0.1 \nabla w_1^2) + 0.1 \nabla w_3^2$$

$$= 0.9 \cdot 0.9^3 \nabla w_3^2 + 0.1 \cdot 0.9^2 \nabla w_1^2 + 0.1 \cdot 0.9 \nabla w_2^2 + 0.1 \nabla w_3^2$$

$$= 0.1 \sum_{i=0}^3 (0.9)^i \nabla w_i^2$$

$\therefore v_t = (1-\beta) \sum_{i=0}^{t-1} \beta^i (\nabla w_i)^2$ , we see that as  $t$  increases, the sum is finite and  $i=0$  so the gradients decays to 0. and the contribution from the early gradients becomes negligible.

In steep areas  $\nabla w$ 's are large hence  $v_t$  also becomes large and the learning rate reduces. On the hand, In gentle areas,  $\nabla w$ 's are small and  $v_t$  decays hence the learning rate increases.

so RMSprop retains the advantages of AdaGrad for sparse and dense features but also has better convergence behaviour because the learning rate can increase or decrease.

One problem RMSprop faces is that it can oscillate around the solution. Near the minima,  $\nabla w$  are small hence making  $v_t$  decay to 0.

This makes the effective learning rate constant  $\frac{\eta}{\sqrt{v_t}}$ , which is susceptible to oscillation. So RMSprop depends on the initialization of these parameters.

⇒ To reduce this effect we can adapt the numerator (as well) this allows us to avoid the dependency on the initializations. This algorithm is called AdaDelta. The algorithm is as follows:-

$$v_t = \beta v_{t-1} + (1-\beta)(\nabla w_t)^2$$

$$\Delta w_t = -\sqrt{\frac{v_{t-1}}{v_t}} \nabla w_t$$

$$w_{t+1} = w_t + \Delta w_t$$

$$v_t = \beta v_{t-1} + (1-\beta)(\Delta w_t)^2$$

Notice that  $\Delta w$  depends on current  $v_t$  and previous  $v_{t-1}$ .

i) In steep areas  $v_t$  is more updated and increases. But  $v_{t-1}$  is not yet updated, so in  $\Delta w_t$ , numerator is smaller than the denominator, hence the learning rate is smaller.  $v_t < v_{t-1}$

2) In ~~flat~~ flat regions,  $v_t$  is more update [and] reduces, but  $v_{t-1}$  is not yet updated; so in  $\Delta w_t$ , the denominator is small, the learning rate becomes larger. ( $v_t \neq v_t$ )

$\Rightarrow$  AdaDelta adapts the learning rate more smoothly than RMSProp algorithm.

Adam (Adaptive Moments): Very similar to the RMSProp algorithm, but introduces momentum in the calculation. So instead of ~~directly~~ using  $\nabla w_t$  in the gradient update, we update it using the momentum term to stabilize the update. and repeat from  $(\beta_1)^t$  for  $m_t$

$m_t = \beta_1 m_{t-1} + (1-\beta_1) \nabla w_t$  where  $\beta_1$  is a constant for scaling/bias removal. for  $\hat{m}_t = \frac{m_t}{(1-\beta_1)^t}$  and  $\hat{v}_t = \frac{v_t}{(1-\beta_2)^t}$  (explained later)

$v_t = \beta_2 v_{t-1} + (1-\beta_2) (\nabla w_t)^2$ , then correcting the bias:  $\hat{v}_t = \frac{v_t}{(1-\beta_2)^t}$

$$\hat{w}_t = w_t - \frac{\eta}{\sqrt{\hat{v}_t}} = w_t - \eta \cdot \frac{\nabla w_t}{\sqrt{\hat{m}_t + \epsilon}}$$

Update rule:  $w_{t+1} = w_t - \eta \cdot \frac{\nabla w_t}{\sqrt{\hat{m}_t + \epsilon}}$

In this momentum, we give more weightage to overall previous gradients in the form of  $\beta_1 m_{t-1}$ , the gradient gradient is given slightly less weightage since it gets multiplied by  $(1-\beta_1)$ . It is like a moving exponential average of  $\nabla w_t$ 's. We can't calculate  $E[\nabla w]$  as it will be expensive, but what about

$$E[m_t] = E\left[(1-\beta_1) \sum_{i=0}^{t-1} \beta_1^i \nabla w_i\right]$$

$$\Rightarrow E[m_t] = (1-\beta_1) \sum_{i=0}^{t-1} \beta_i^t E[\nabla w_i]$$

$E[m_t] \approx E[\nabla w]$ , therefore  $\frac{m_t}{1-\beta_1}$  gives an unbiased estimate of  $\nabla w$ .

$\Rightarrow$  A lack of initialization bias correction would lead to initial steps that are much larger. Bias correction also ensures that initial learning rate is not very huge.

Max Prop :- If we closely revisit the formula for RMSPROP:-

$$v_t = \beta v_{t-1} + (1-\beta)(\nabla w_t)^2, \quad w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \nabla w_t$$

Although it achieves the requirement of adjusting the learning rate based on the gradient. But is still aggressive and that makes it susceptible to noise. Take for example a sparse feature in a Mini-Batch GD and the gradient is considerable. Now a batch comes which was noisy and no data point had this feature, so gradient will be zero, in the current rule, the zero gradient which is noisy will affect  $v_t$  and reduce it and hence increase the learning rate. This is undesirable.

$\Rightarrow$  We change the way  $v_t$  is calculated to make it less susceptible to noise. That is done in the max prop instead of using the

$L^2$  norm (squares) we can use the max norm, so  
the update formula becomes:-

$$v_t = \max \{ \beta v_{t-1}, |\nabla w_t| \}$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_t + \epsilon}} \nabla w_t$$

square root is removed as max norm is used now!

Thinking about the same situation, zero gradient will never be picked, but still the  $v_t$  decays as  $\beta$  is multiplied, so if the behaviour is consistent,  $v_t$  will decay and the original behaviour is still retained.

~~But in the case of Adam :- start from scratch~~

Introducing Max norm in Adam :- Adamax

$$m_t = \beta m_{t-1} + (1-\beta) |\nabla w_t|$$

$$v_t = \max \{ \beta v_{t-1}, |\nabla w_t|^2 \}$$

notice that  $\hat{v}_t$  is not needed as max norm does not drastically increase the learning rate.

NADAM

Introducing Nesterov Acceleration in Adam :-

Slightly different notation (did not understand)

$$m_{t+1} = \beta m_t + (1-\beta) \nabla w_t$$

$$\hat{m}_{t+1} = \frac{m_{t+1}}{1-\beta^{t+1}}$$

$$v_{t+1} = (\beta_2 v_t + (1-\beta_2) (\nabla w_t)^2)^{1/2} = \frac{v_{t+1}}{1-\beta_2^{t+1}}$$

$$w_{t+1} = w_t - \frac{\eta}{\sqrt{v_{t+1} + \epsilon}} \left( \beta_1 \hat{m}_{t+1} + (1-\beta_1) \nabla w_t \right)$$

## Learning rate Schemes

<u>Based on epochs</u>	<u>Base Validation</u>	<u>Based on gradients</u>
1) Step Decay	1) Line Search	1) AdaGrad
2) Exponential Decay	2) Log Search	2) RMSProp
3) Cyclical	3) Adam	3) Adadelta
4) Cosine Annealing	4) Adam	4) AdamW
		5) Adamax
		6) Nadam
		7) AMSGrad
		8) Adam

Cyclical Learning Rate :- At saddle points, the gradients become small, and it becomes difficult to move out from those areas.

If we allow the learning rate to increase after some iterations irrespective of the gradients, then there is a chance to escape the saddle points. Adaptive learning rate helps in escaping saddle points but they are computationally expensive.

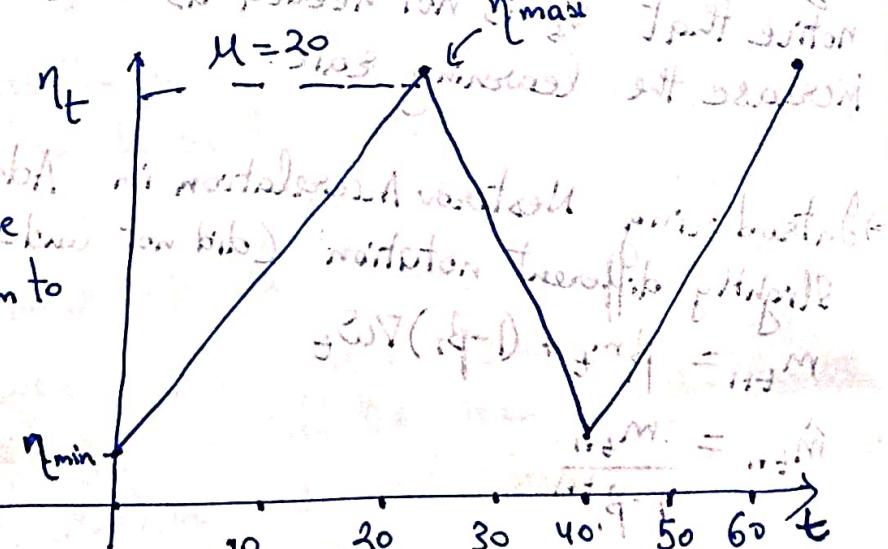
We can vary learning rate cyclically as a simple strategy.

i) Triangular :-

After every 20 steps, the learning rate reaches from minimum to maximum.

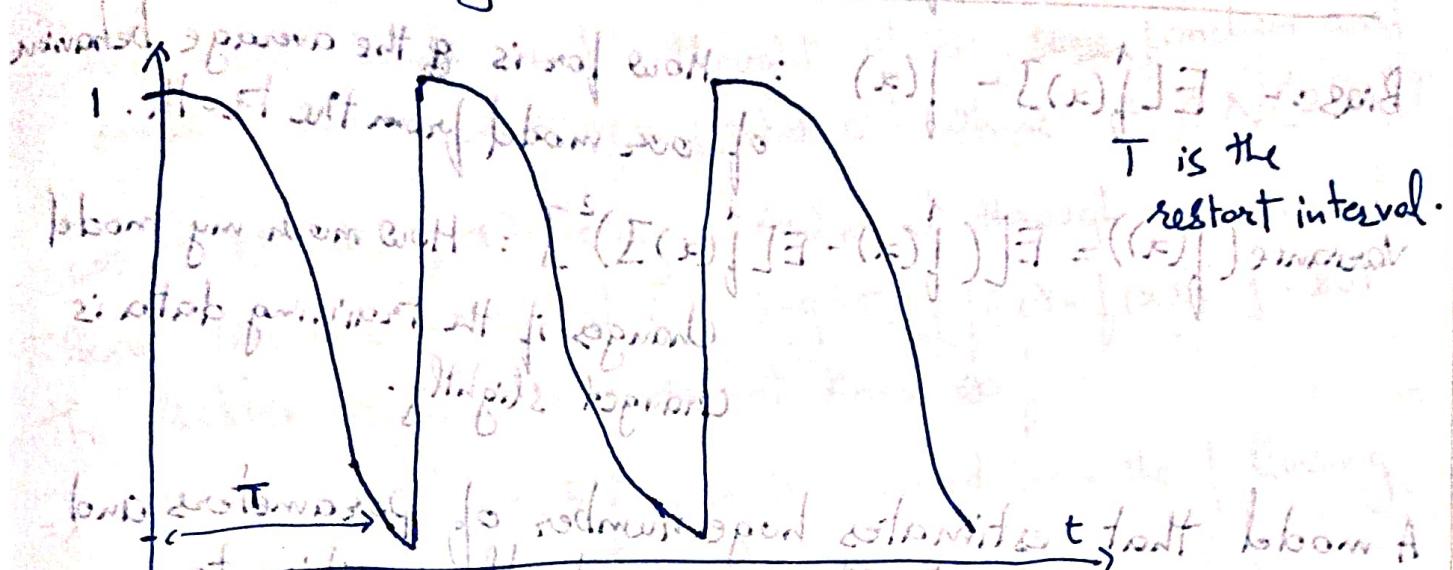
Formula is given by:-

$$\eta_t = \eta_{\min} + (\eta_{\max} - \eta_{\min}) \cdot \max(0, (1 - |\frac{t}{20} - (2L + 1)|) + 1)$$



$$\left( \frac{\eta_{\max}(\eta_{\max} - \eta_{\min})}{20} \right) \frac{t}{20} + \eta_{\min} = \eta_t$$

## ii) Cosine Annealing (Warm-restart)



$L(x) = L(0) + [L(0) - L(T)] \cdot \cos(\frac{\pi t}{T})$

$L(0) = \min_{t=0} L(t)$

$L(T) = \max_{t=T} L(t)$

But due to ~~the~~ abrupt increases in the learning rate, it may need more epochs. To avoid overshoot near the minima. In that case, we could use techniques such as early stopping, to roll back to the minimum.

Warm Start (found to be helpful in transformed architectures).

Typically, we set the initial learning rate to a high value and then decay it. On the contrary, using a low initial learning rate helps the model to warm up and converge better. This is called warm start.

This is because, initially the weights are random and we allow the model to make small steps in early gradient directions. Then after certain Number of iterations, the usual

exponential decay strategy is applied. Number of warm-up steps until is

bounds can be a choice.

## Bias and Variance of a Model

Bias :-  $E[\hat{f}(x)] - f(x)$  : How far is the average behavior of our model from the truth.

Variance ( $\hat{f}(x)$ ) =  $E[(\hat{f}(x) - E[\hat{f}(x)])^2]$  : How much my model changes if the training data is changed slightly.

A model that estimates huge number of parameters and capable to fit complex curves, is highly sensitive to training data. If this model is trained on a slightly different sample of training data, it will learn significantly different learned parameters. Thus it has high variance. Since this model can fit the training data, it leads to low error and hence low bias as it has the capability to fit the truth closely. But sometimes, if the data has noise, then this model gets severely affected because even though the training error will be small as always, but, on unseen data, it performs poorly. This is called overfitting. If still has low bias (usually) as the model is capable to fit the truth if we ignore the noise. So overfitting is a consequence of high variance and noise.

$\Rightarrow$  Similarly, a model that estimates very few number of parameters. Thus even if the training data sample is slightly changed, the learned curve remains similar. Thus low variance. If the model is too simplistic, then it may fail to capture the structure in the data, thus leads to high bias as the model does not have the capability to fit the truth irrespective of training data. This leads to high training error and is called underfitting.

## Estimating error from test data

the true labels can be thought of as  $f(x)$  function and some noise. So,  $y = f(x) + \epsilon$  where  $\epsilon \sim N(0, \sigma^2)$

We estimate  $f(x)$  using  $\hat{f}(x)$  or  $\hat{y}$  through our model. Hence we are interested in using  $E[(\hat{f}(x) - f(x))^2]$ . But the problem is that we do not know  $f$ .

$\Rightarrow E[(\hat{f}(x) - f(x))^2]$  can be derived from the following formula: it is  $E[(\hat{y} - y)^2] = E[(\hat{f}(x) - f(x) - \epsilon)^2]$

$$= E[(\hat{f}(x) - f(x))^2] + E[\epsilon^2] - 2E[\epsilon(\hat{f}(x) - f(x))]$$

$$\Rightarrow E[(\hat{f}(x) - f(x))^2] = E[(\hat{y} - y)^2] + 2E[\epsilon(\hat{f}(x) - f(x))]$$

$$= \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + 2E[\epsilon(\hat{f}(x) - f(x))].$$

For 2 independent variables  $X$  and  $Y$ , and one having expectation  $E[XY] = E[X] \times E[Y] = 0$  (i.e. either  $E[X]$  or  $E[Y]$  is 0)

Since in test data, none of the points participated in estimating  $\hat{f}(x)$ ,  $(\hat{y}_i - \hat{f}(x))$  is independent from  $(\hat{y}_i - f(x))$ . Hence.

$\epsilon$  is independent from  $(\hat{f}(x) - f(x))$ .

$$\text{So } E[(\hat{f}(x) - f(x))^2] = \frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 - \bar{\epsilon}^2 \quad (\bar{\epsilon} \text{ depends on the noise.})$$

If the data is sufficiently clean, then this is a good estimate of the true error (LHS).

If instead we used training data for estimation then  $E[\epsilon_i | \mathcal{D}]$  is not independent from  $f(x) - \hat{f}(x)$ . Hence  $E[\epsilon_i^2 | \mathcal{D}] \neq 0$ . The true error LHS is larger than the error.

So here we want  $\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(y_i - \hat{y}_i)^2]$ .

So using the empirical training error does not give a true representation of the actual error. But training is based on minimizing this error.

Stein's Lemma:  $\mathbb{E}\left[\sum_{i=1}^n \epsilon_i (\hat{f}(x_i) - f(x_i))\right] = \frac{1}{n} \sum_{i=1}^n \frac{\partial f(x_i)}{\partial y_i}$

If  $\frac{\partial f(x_i)}{\partial y_i}$  is large when there is high variance in the model, changing  $y_i$  brings large change in the learned model.

So LHS of Stein's Lemma is similar to  $\mathbb{E}[E(\hat{f}(x) - f(x))]$

So we can say that true error is a function of model complexity.

$\mathbb{E}[\hat{f}(x) - f(x)] = \text{empirical train error} - \text{small constant} + \Omega(\text{model complexity})$

The more complex the model, the true error may increase.

The complex the model, the modified objective function

so it is better to minimize a modified objective function

that also incorporates model complexity into it.

$\min_{\theta} L(\theta) = \text{train error} + \Omega(\theta)$ . But instead of calculating

$\Omega(\theta)$  given by Stein's lemma, we can calculate it using an approximation that captures model complexity. This strategy is called regularization.

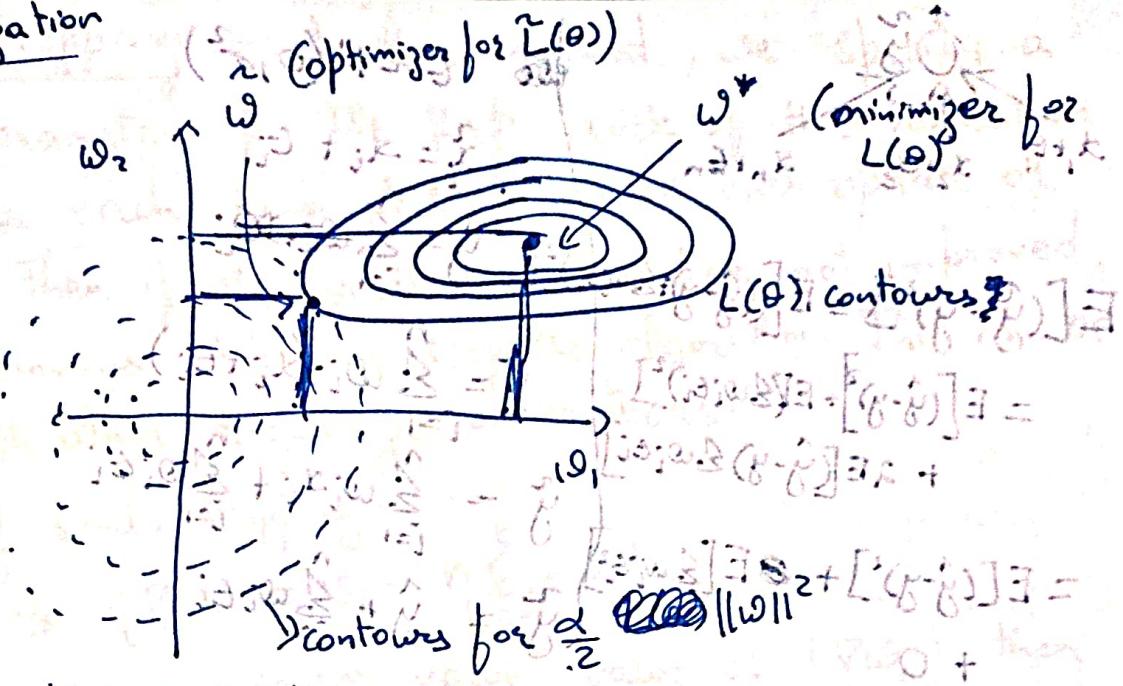
We can calculate it using an approximation that captures model complexity. This strategy is called regularization.

$L(\theta) = \mathbb{E}[(y - \hat{y})^2] = \mathbb{E}[(y - f(x))^2]$ , where

$\hat{y} = f(x) + \epsilon$ . Here  $\epsilon$  is a random variable with mean zero and variance  $\sigma^2$ .

$\mathbb{E}[(y - f(x))^2] = \mathbb{E}[(f(x) + \epsilon - f(x))^2] = \mathbb{E}[\epsilon^2] = \sigma^2$ .

## i) L<sub>2</sub> regularization



the eigen value corresponding to  $w_1$  will be smaller. Hence the ellipse is elongated in  $w_1$  axis. Moreover, there is more shrinkage in  $w_1$  as compared to  $w_2$ . This can also be inferred through the contours of  $L(\theta)$ , because in the axis of  $w_1$ , the contours are further apart, so derivative is small. So loss is less sensitive to  $w_1$  than  $w_2$ . So  $w_1$  is shrunk more.

ii) Data Augmentation:- we can introduce new training data points by slightly perturbing the already present data points and then keeping their labels same. Example, shifting or rotating an image to create a new image.

iii) Injecting noise at inputs:- at each epoch, corrupt some part of the input, thus it becomes difficult for the model to overfit on the training data as at each epoch the training data slightly changes. Can be thought of data augmentation at each epoch. For an input output network, without any hidden layers, the noise injection at input is equivalent to L<sub>2</sub> regularization.

$$\begin{aligned}
 & \text{Diagram showing } \hat{y} = \sum_{i=1}^n w_i x_i + b \\
 & \text{Inputs: } x_1, x_2, \dots, x_n \\
 & \text{Weights: } w_1, w_2, \dots, w_n \\
 & \text{Bias: } b \\
 & \text{Error: } (y - \hat{y})^2 \\
 & \text{Expected value: } E[(y - \hat{y})^2] = E[(\hat{y} - y)^2] \\
 & = E[(\hat{y} - y)(\hat{y} - y)] = E[(\hat{y} - y) + E[\sum w_i e_i]]^2 \\
 & = E[(\hat{y} - y)^2] + E[\sum w_i e_i]^2 \\
 & + 2E[(\hat{y} - y)\sum w_i e_i] \\
 & = E[(\hat{y} - y)^2] + E[\sum w_i^2 e_i^2] \\
 & + 0
 \end{aligned}$$

This is because  $E[(\sum w_i e_i)^2] = E[\sum w_i^2 E[e_i^2]]$   
 since  $e_i$  is independent from  $e_j$  and  $E[e_i e_j] = E[e_i] E[e_j] = 0$   
 So, all interacting terms in LHS will be 0.  
 Hence only square term will survive.  
 Also  $e_i$  is independent from  $(\hat{y} - y)$  so  $E[(\hat{y} - y)e_i] = 0$   
 $= E[(\hat{y} - y)^2] + E[\sum w_i^2 E[e_i^2]]$

$$\begin{aligned}
 & \text{So } E[(\hat{y} - y)^2] = L(\theta) + E[\sum w_i^2 e_i^2] \\
 & \boxed{E[(\hat{y} - y)^2] = L(\theta) + \sum w_i^2}
 \end{aligned}$$

v.) Injecting noise at the outputs:— Instead of fitting on the true labels, we perturb the true labels by adding some noise.

This has a similar effect to L2 regularization.

This injects extra steps to the loss function, so it's not smooth.

v) Early Stopping: In this method, we keep specifying a patience parameter  $p$ . We keep track of the validation error after each step. A step can be each update or each epoch. Then if the validation error has not improved after  $p$  continuous steps we stop the algorithm. This can be combined with other regularization methods.

Analysis of Early Stopping:-  
 $\omega_{t+1} = \omega_t - \eta \nabla \omega_t$  is a good approximation based on (s)  $\omega_{t+1} = \omega_t - \eta \sum_{i=0}^t \nabla \omega_i$  is a good approximation based on (s)  $\omega_{t+1}$  is stable based on value of  $|\nabla \omega_i|$ , then let  $T$  be the maximum value (absolute)

of  $t$  such that  $|\nabla \omega_i| \leq \eta T$  this step controls the  $|\omega_{t+1} - \omega_0| \leq \eta T$  this step controls the  $|\omega_{t+1} - \omega_0|$ . If we limit  $t$ , in some upper bound of  $|\omega_{t+1} - \omega_0|$ . If we limit  $t$ , in some sense we are limiting  $|\omega_{t+1} - \omega_0|$  and if  $\omega_0 = 0$ , then this acts as a regularizer on the weights.

vi) Ensemble Methods:- bagging does not help if the ensembles are correlated as they will make similar mistakes. Bagging is most effective when ensemble models are decorrelated as mistake of one

rectified by the other model

2 options available in bagging are:-  
i) Train different neural networks having different architecture  
ii) Train multiple instances of the same network using different training samples.

Both these options are computationally expensive.

vii) Dropout:- In this method we train a single instance of neural network, but at each step (update or epoch), we drop some neurons by temporarily removing their incoming and outgoing edges. Neurons to drop are decided randomly.

Sampled where some probability is associated with each neuron. For example, each neuron has  $p=0.8$  probability of being retained. The total number of possible neural networks that can be made, are exponential in the number of total nodes. The full algorithm is given below:-

- 1) Initialize the network by choosing which neurons are active and participate in the computation.
- 2) Do forward propagation using the weights from previous iterations.
- 3) Do backward propagation and update the weights that are active.
- 4) Repeat the procedure. We can decide if we want to keep one network or all. update or 1 epoch. After each step, a new or in rare cases a previous configuration is made active and weights are updated.

Each network is trained twice rarely, but each weight may be updated many times in different configurations.

At test time, the output is calculated from the whole network but the output of each neuron can be multiplied by the probability it was retained with in each sample, which is  $p$ , or by the fraction of times it was active during training.

Dropout essentially applies noise to the hidden units. It prevents units from co-adapting. A hidden unit cannot rely too much on other units as they may get dropped out any time. It allows the model to discover new features and makes it more robust.

Unsupervised Pre-Training: Instead of initializing the parameters randomly, is there a way to initialize them? Consider the first hidden layer, having some  $n$  neurons. If we ignore the rest of the network and set up a new network such that the task of this network is to reconstruct the input layer with the least error possible with only the 1<sup>st</sup> hidden layer consisting of the neurons from 1<sup>st</sup> hidden layer of the original network.

So the objective is such that:-

$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m (\hat{x}_{ij} - x_{ij})^2$$

Each input is  $R^n$ , and there are  $n$  such inputs. The reconstructed output is also  $R^n$ . This is VPCA, the objective is same. Whatever weights are learned as  $W_1$  (between input to 1<sup>st</sup> hidden layer) can be used as initialization in the original network.

For the weights of the second layer ( $W_2$ ) we calculate the output of each neuron in the first layer for each data point using the weights in the first step. So output of the 1<sup>st</sup> hidden layer will be  $R^m$ , and there are  $n$  such inputs.

$$\text{The objective will now be: } \min \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m (\hat{h}_{ij} - h_{ij})^2$$

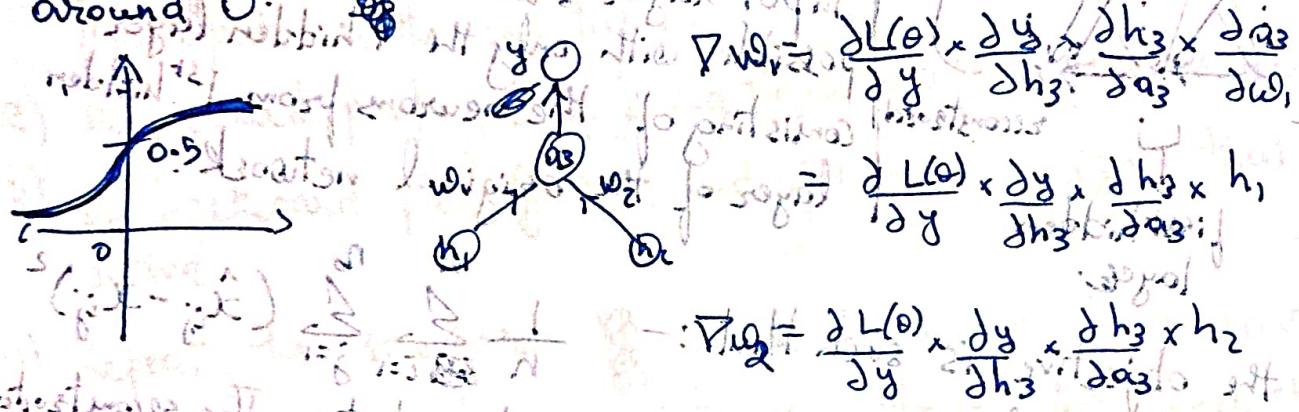
$h_{ij}$  is the output from the  $j^{th}$  neuron of the  $i^{th}$  point of  $n$  vectors.  $\hat{h}_{ij}$  is the output from the  $j^{th}$  neuron of the  $i^{th}$  point of  $n$  vectors. From this we obtain  $W_2$ .

We do this to obtain  $W_1$  and then initialize these as the weights for the original network.

Activation Functions

- 1) Linear: If we have a deep linear neural network, then it faces difficulty in creating approximate all the functions. Therefore some non-linearity is needed.

- 2) Sigmoid: Sigmoid function compresses all its domain into  $[0, 1]$ . The output of sigmoid is positive and it is not centered around 0.



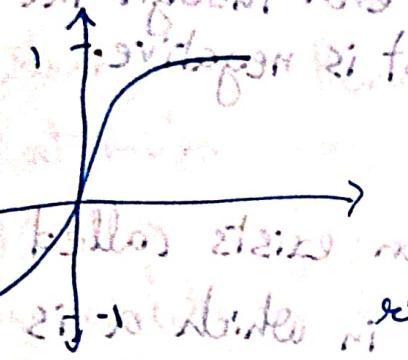
Since neurons are sigmoid,  $h_1$  and  $h_2$  are positive. Besides  $h_1$  and  $h_2$ ,  $D_w_1$  and  $D_w_2$  are same. Hence both  $D_w_1$  and  $D_w_2$  will have the same sign. Second and fourth quadrant gradients cannot be easily calculated as signs of the components

$D_w_2$  (Allowed) are opposite. It is mitigated by centering  $w_2$  removes this problem.

One more problem with sigmoid neuron is that near large absolute values of input, the neuron outputs 0 or 1, but the gradient vanishes as the rate of change is very small. So incorrect initializations of weights makes the neuron saturate (0, 1). And after saturation, since the gradient is small, the updates happen very small which keeps it saturated in the next iterations.

Third problem is that sigmoid involves sigmoid calculations, which is also expensive.

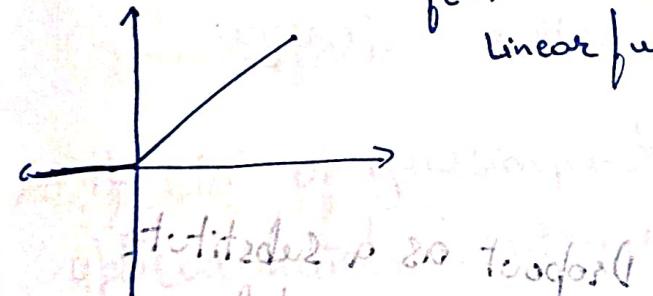
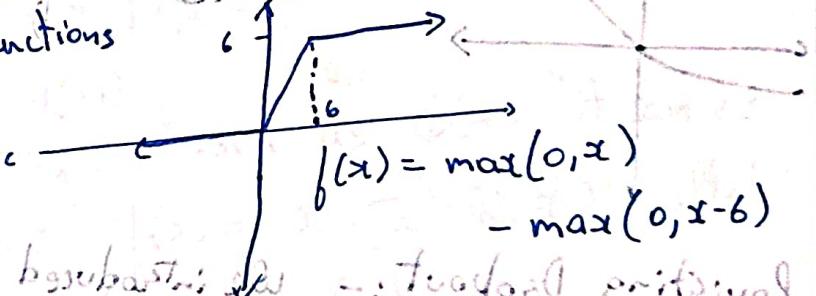
3) Tanh function:  $\tanh(x)$ . without  $b$  it's called hyperbolic tangent function.  
 $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



compresses all its input to the range  $[-1, 1]$ . It is zero centred so it removes the opposite sign gradient. But the problem of vanishing gradient still exists. And it also involves exponential computations, so it is also relatively expensive.

4) ReLU (Rectified Linear Unit)

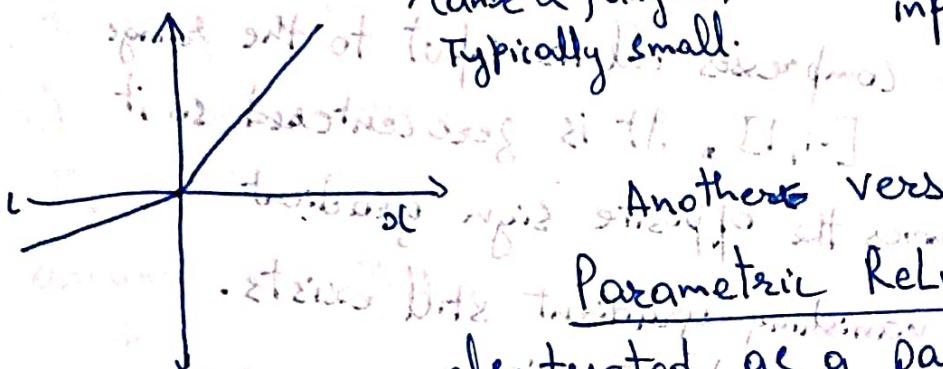
$f(x) = \max(0, x)$  We can recover piecewise linear functions

Advantages of ReLU:  
 Does not saturate in the positive region, computationally efficient, since it converges much faster than sigmoid/tanh.  
 $\frac{\partial \text{ReLU}(x)}{\partial x} = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$

If bias of some ReLU is set to a large negative value even  $b < 0$ , then the output as well as gradient of this neuron will be 0. Hence the vanishing gradient problem still exists for one half of input.

Leaky Relu : No Saturation. (Without : without shift) Some gradient value still exists even though the input is negative.

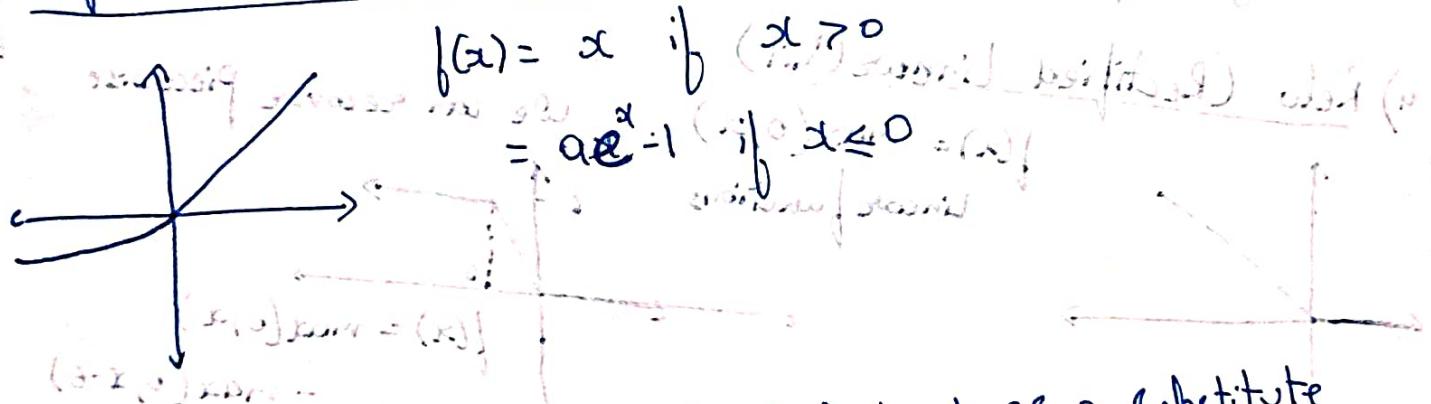


Another version exists called Parametric Relu in which  $\alpha$  is

also treated as a parameter, which gets

updated during backpropagation.

## Exponential Linear Unit



Revisiting Dropout :- We introduced Dropout as a substitute to bagging. But in bagging, each individual model is overfit on its corresponding bagged training data. But in dropout, each subnetwork configuration hardly fits on the data. To achieve the bagging effect with dropout, we introduce MaxOut neuron.  $\Rightarrow \max \{w_1 x + b_1, \dots, w_n x + b_n\}$ . Only the neuron that contributes the most gets the update, as for others, gradient is 0.

Dropsout can also be written as  $O(x) = \begin{cases} 1 \cdot x & \text{if } \text{rand}(0,1) < p \\ 0 \cdot x & \text{if } \text{rand}(0,1) \geq p \end{cases}$ . Output of a neuron is as it is ( $x$ ) if it got selected, the probability of that is  $p$ . And output is 0 otherwise.

It makes sense to combine and represent more compactly:-  $f(x) = m \cdot x$

where  $m \sim \text{Bernoulli}(\phi(x))$  where  $\phi(x)$  is a function that maps  $-\infty, -\infty$  between  $[0, 1]$ : ~~For dropout~~ ~~most learned~~

$\phi(x) = p$  works, as each neuron will be selected with probability of  $p$ .

Alternatively it can be a cumulative distribution function of a Gaussian variable or even logistic function.

$$E[f(x)] = E[m \cdot x] = x \cdot E[m]$$

$$\text{Therefore } = \mathbb{E}[x \cdot \phi(x) + 0 \cdot (1 - \phi(x))] \text{ using principle}$$

$$(Since: E[f(x)] = \mathbb{E}[x \cdot \phi(x)] \rightarrow P(X \leq x))$$

if CDF of gaussian function is chosen as  $\phi(x)$ , it can be approximately written as  $f(x) = x \cdot \sigma(1.702 \cdot x)$  called GelU.

SELU (Scaled Exponential Linear Unit): Linear Unit

with output scaled such that mean is 0 and unit variance.  $f(x) = \begin{cases} \alpha e^{\alpha x} & \text{if } x > 0 \\ \alpha e^{-\alpha x} & \text{if } x \leq 0 \end{cases}$   $\alpha \approx 1.677363$

Studies found that functions of the form  $f(x) = x^\beta (\beta x)$  are good. Function of such forms where  $\beta$  is a learnable parameter is called SWISH. if  $\beta = 1.702$  it is GelU.

when  $f(x) = x \phi(x)$ , it is called as Sigmoid weighted linear unit.

$$x \cdot \phi(x) = (x)_\phi \rightarrow \text{Handwritten}$$

where  $\phi(x)$  makes  $(x)_\phi$  differentiable or make

Convolution operation: If  $I$  is inserted in  $\phi$  form to  $K$  then  $I$  is  $\begin{bmatrix} m \\ n \end{bmatrix}$  and  $K$  is  $\begin{bmatrix} m \\ n \end{bmatrix}$

$$S_{ij} = (I * K)_{ij} = \sum_{a=-\frac{m}{2}}^{\frac{m}{2}} \sum_{b=-\frac{n}{2}}^{\frac{n}{2}} I_{i-a, j-b} \cdot K_{\frac{m}{2}+a, \frac{n}{2}+b} \cdot \phi \text{ for } g_t \text{ is derivative}$$

position, position of  $I$  with  $K$  is  $\phi$  so  $I$  has to be  $3D$  as  $K$  is  $3D$ . In case of  $3D$  filter is referred to as a volume. In case of  $2D$  filters, the volume is slid over the  $3D$  image same as  $2D$  case, and the convolution result is  $2D$ .

Defining Quantities:

$\Rightarrow$  Width ( $W_1$ ), Height ( $H_1$ ) and Depth ( $D_1$ ) of the original input

$\Rightarrow$  The stride  $S$  (increment in the slide of the kernel)

$\Rightarrow$  The number of filters  $K$  is,  $K$  is misspelled as  $F$ .

$\Rightarrow$  The spatial extent ( $F$ ) of each filter (the depth of each filter is same as the depth of each input)

$\Rightarrow$  The output is  $W_2 \times H_2 \times D_2$ .

Usually the extent of the filter is an odd integer as center of the matrix is properly defined.

If no padding is applied, the convolution reduces the size of the image as, the pixels in the border of the input (convolution is not defined).

The formula for output dimensions is given as:-

$$W_2 = W_1 - (F - 1) = W_2 - F + 1 \text{ where } F \text{ is filter size}$$

$$H_2 = H_1 - F + 1 \text{ if } F \text{ is filter size}$$

This is because if  $F = 2a+1$  for  $a \geq 0$ ,  
 $a$  columns are lost in the left border and  $a$  columns are  
lost in the right border of the input.  $\therefore$  total lost columns =  $2a = F-1$   
 $\therefore$  remaining =  $W_1 - (F-1)$

Similarly for height, remaining height =  $H_1 - (F-1)$

If we add padding of  $P$  pixels around the image, then  
the formula is updated as follows:

$$W_2 = W_1 - F + 1 + 2P$$

$$H_2 = H_1 - F + 1 + 2P$$

Generally we do  $s=1$  which is to look at all pixels.  
 $s=2$  means that place kernel at every alternate pixel. (in both  
x and y dimension)

$$W_2 = \frac{W_1 - F + 2P}{s} + 1$$

$$H_2 = \frac{H_1 - F + 2P}{s} + 1$$

$\Rightarrow K$  filters will give us  $K=2D$  outputs. Hence  $D_2 = K$

## Convolutional Neural Networks

In traditional Machine learning, we would convert the image  
into a flat vector and then apply the classification algorithm.  
We can also apply edge detection algorithms and kernels.  
before applying the classification algorithms. But these  
kernels are hard-coded before applying the algorithms and is  
treated completely as preprocessing.

We can also learn these kernels and treat them as params  
of the model. Not only this, multiple kernels can be learnt.  
Furthermore the neural network's architecture allows to learn  
multiple layers of kernels to produce the final output.

## Difference between a CNN and a normal feed-forward Network

In a feed-forward network, each neuron is connected to each neuron in the previous network. It is a dense network and all weights are independently treated.

But in a CNN:-

Suppose there is an image of  $4 \times 4$  given as input to the network, and a  $2 \times 2$  kernel is trying to be learnt for the first hidden layer.

i) The input image is flattened :-

Hence input layer has 16 neurons.

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

$$\Rightarrow [1|2| \dots |15|16]$$

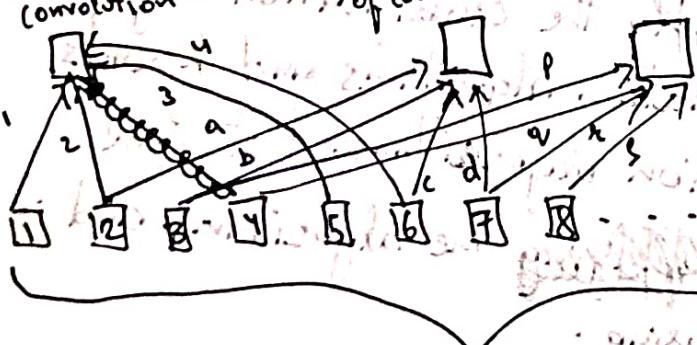
ii) The kernel is flattened :-

1	2
3	4

1	2	3	4
---	---	---	---

iii) The output image is calculated and flattened as well. So each pixel in the output of the convolution is a neuron in the next layer.

We can make the following observations:-



i) The output layer is sparsely connected with the previous layer. This is because the first pixel of convolution output is calculated using only pixels 1, 2, 5 and 6 of input.

Similarly for all pixels. So the connectivity is sparse.

ii) Since a single kernel is applied here,  $1=a=f$ ,  $2=b=g$ ,  $3=c=r$ ,  $4=d=s$ .

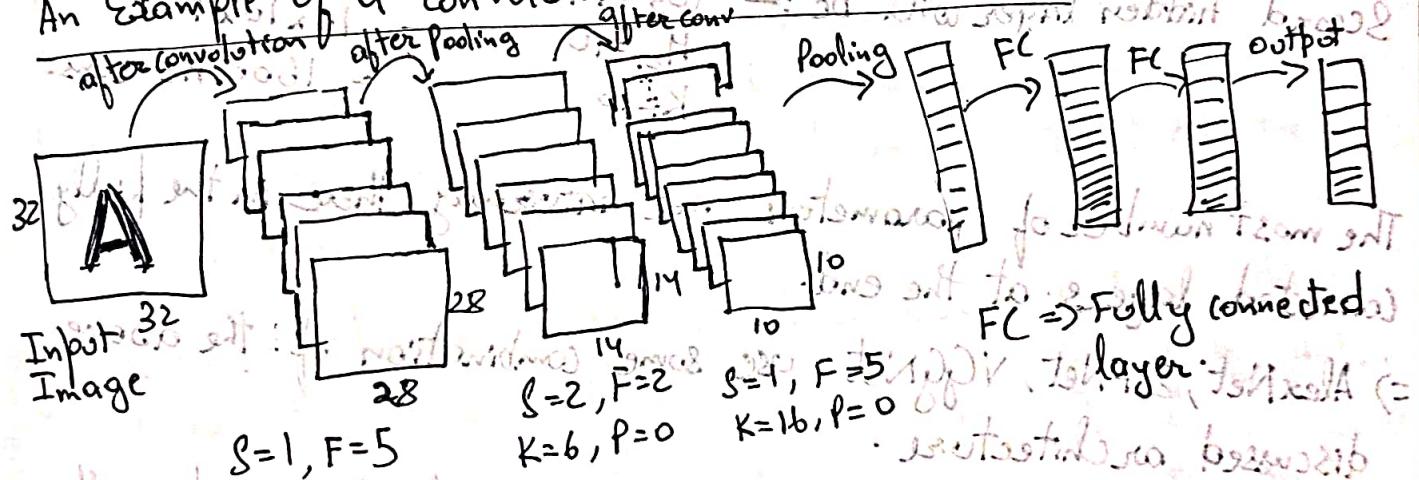
There is weight sharing in this neural network. More kernels can be applied but that also increases the number of output neurons (next layer convolution output, addition in depth) so weight sharing is still there.

Due to this weight sharing:- i) we look at the neighborhood of each pixel.

ii) We don't really lose information as in the next layers, there will be interaction between pixels which depend on the input pixels. So indirect interaction is there.

=> After performing convolution the result can be passed through an activation function.

An Example of a Convolutional Neural Network:-



The first layer consists of 6 such filters, so there are  $6 \times 25 = 150$  parameters. Note that this means that there will be 150 "unique" / independent connections from the input layer to the first hidden layer. Moreover there are  $28 \times 28 \times 6$  neurons in the first hidden layer. This is decided by the formula discussed before.

The Layer after first hidden layer is the max pooling layer. In this layer is passed through a max filter independently to obtain the corresponding matrices. Notice that a stride of 2 is used in the example, dimension rules apply in the same manner.

In max filter, just the max is selected within the neighbourhood defined by the kernel.

Max pooling is done to reduce the size of the layers and hence reduce the number of parameters, preventing overfitting.

Also note that the max pooling layer learns no parameters, it is a fixed process involving no learnt parameters.

=> The next hidden layer (second), is again a convolutional layer with filter size of  $5 \times 5 \times 6$ . There are 16 such filters. Hence the number of parameters are 2400. The number of neurons in the second hidden layer will be :-



$$W_2 = 14 - 5 + 1 = 10$$
$$H_2 = 10$$
$$K = 16$$
$$= 10 \times 10 \times 16 = 1600 \text{ neurons.}$$

The most number of parameters are naturally there in the fully connected layers at the end.

=> AlexNet, ZFNet, VGGNet use some combination of the above discussed architecture.

Notice that in each layer, if we are using multiple kernels, all those kernels are fixed same, both in regard of type of Kernel and size of Kernel. We can combine these options and the concatenate them to produce the output. but this introduces a problem, if a filter is  $F \times F \times D$ , then each entry in the output requires  $F \times F \times D$  computations. This makes it difficult to apply multiple kernels. We can use  $1 \times 1$  convolutions to reduce the depth of the input.

$H \rightarrow$  The  $(1 \times 1) \times D$  convolution operation applies across the depth of the input.



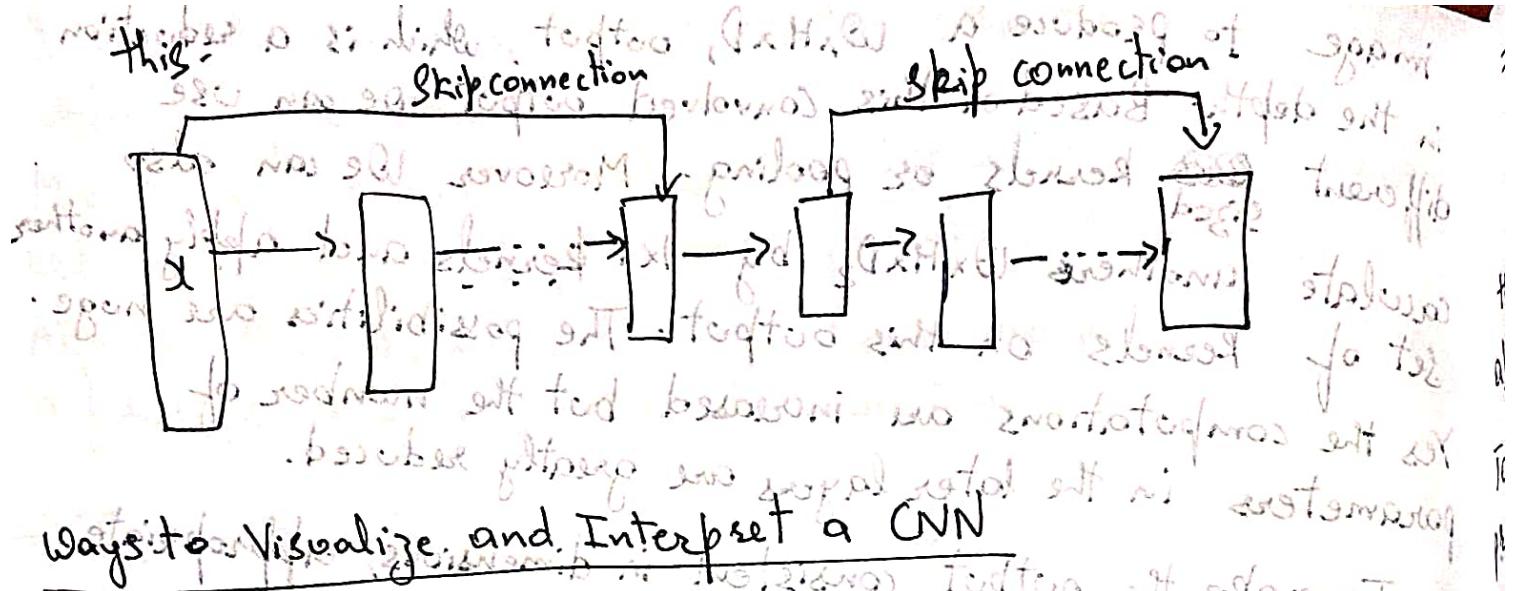
=> We can apply  $D$  such  $1 \times 1$  kernels to the input

image to produce a  $W \times H \times D$ , output, which is a reduction in the depth. Based on this convolved output we can use different sized kernels or pooling. Moreover we can also calculate another  $W \times H \times D_2$  by  $1 \times 1$  kernels and apply another set of kernels on this output. The possibilities are huge. Yes the computations are increased but the number of parameters in the later layers are greatly reduced.

To make the output consistent in dimensions, appropriate padding is added. All these computations are done in parallel. This whole unit is called an inception module. The unit consists of  $1 \times 1$  convolutions and pooling and other size convolutions. The output is concatenated and treated as the output of the inception module as a single layer.

One more change can be to do average pooling in the last layer before the fully connected layer. So if it has  $E \times F \times D$  output, then average pooling reduces it to  $D$  outputs, applied by taking average of each component matrix. This reduces the number of parameters. All these improvements were used in the GoogleNet.

ResNet: Suppose we have a shallow Network that performs well... but if we make it deeper, we expect it to perform at least as better as the shallow network. But it does not happen in practice. To improve the performance of the deep network, we can use skip connections. Skip connections supply a copy of input to further layers. We can introduce skip connections after every few layers. In fact, ResNet did



## Ways to Visualize and Interpret a CNN

- => For each neuron or set of neurons, we can feed in images to this CNN and identify the images which cause these neurons to fire.
- => We can then trace back to the pattern in the image which causes these neurons to fire. This will allow us to understand which neuron(s) are fired at what input.
- => One way to do this is to plot the kernels/filters as images and interpret them as patterns. The filters detect the patterns observed in their image plot. The problem with this approach is that the results are interpretable only for the first layer. Because after that, the filters are applied on the input of previous layers.
- => Second approach is through occlusion experiments. In this method, we occlude (gray out) different patches in the image and see the effect on the predicted probability of the correct class. The areas that affect the probability the most are the most important in the image.

=> Third approach is to calculate the gradients of a neuron with respect to each input pixel. This will us an image of gradients. (Each image can be seen and wherever the value of gradient is high, more white pixels will appear in the output of this step with respect to that neuron. To calculate these gradients we can use backpropagation. One problem with this method is that gradients can be highly positive or negative in areas of importance, but when visualised through a digital image, these ~~area~~ gradients cancel the effect of each other in these areas and that leads to blunt results and noisy output.

To deal with this problem, we use guided backpropagation. In this method, we choose a neuron of interest and set all other neurons in that layer to 0. Then we input an image to the network and perform a feed forward step. Then we perform the backpropagation step to calculate the gradient matrix w.r.t. to the input. But instead of allowing negative gradient to flow back, we clip negative gradients to 0. This is done to prevent all neurons in all the layers. This prevents the effect discussed earlier and the results become more vivid in the resultant image. Basically we focus only on the positive influences (gradients).

=> One more method to interpret the results of a CNN is to optimize over images. The goal here is to create an image which gets classified as a certain class. So, if this image is passed through a trained convNet it should maximize the probability of the class of interest. Doing this would allow us to visualize what images are classified in which manner.

To do this we can set this as an optimization problem with respect to the pixels of the input image. The reason is we have a need of the argmax ( $-L(\theta)$ )

$$\text{argmax} (\text{Score}(I) - \lambda \cdot \mathcal{R}(I))$$

Now doing with the score, first is the loss for pixel with respect to the softmax. where  $\text{Score}_c(I)$  = score for class before softmax and  $\mathcal{R}(I)$  = some regularizer to ensure that  $I$  looks

real-like and image is shriveling. Now is better with this. In the whole process we keep the parameters of the convolutional neural network fixed. We adjust the image pixels so that the score of the class is maximized.

=> We start with a zero image / randomly initialized image. Now at

=> Set the score vector to be  $[0, 0, \dots, 100, 0, 0]$  better with this.

=> ~~Do a forward pass.~~ Do a forward pass.

=> Compute the gradients  $\frac{\partial L(\theta)}{\partial I_{i,j}}$ , if all is justified then shriveling

now of the image with pixel  $i, j$ . Note that it's propagated

=> Update the pixel  $I_{i,j} = I_{i,j} - \eta \frac{\partial L(\theta)}{\partial I_{i,j}}$ .

=> Repeat the process from forward pass of image until

Instead of limiting this process to only the output, we can focus on any one neuron in any layer to see what images makes that neuron to fire.

To do this

=> Feed an image through the network. (The propagation can stop at the neuron of interest)

=> Set activation  $\sigma$  in layer of interest to all zero, except

for the neuron of interest. This is sigmoid of

=> Backprop to image

=>  $I_{i,j} = I_{i,j} - \eta \frac{\partial L(\theta)}{\partial I_{i,j}}$ , where  $L(\theta)$  can be negative of

the pre-activation of the neuron of interest. It is called

Another similar idea is to create images using embeddings.

So there is a pre-trained CNN, now the objective is to feed this network with an image, not the embeddings of each layer, and then create/reconstruct an image such that the created image's embeddings are as close to the original image's embeddings.

Embeddings are the output of each layer. So the loss function is the squared error between the embeddings of the original image and the created image.

$$L(I) = \|\phi(I) - \phi(I_0)\|^2 + \lambda \|\phi(x)\|_2^6$$

- $\phi(I)$  = vector of embeddings of the created image.
  - $\phi(I_0)$  = vector of embeddings of the original image.
  - $I[i, j] = I[i, j] - \eta L(I)$
- $\Rightarrow$  This can be done for embeddings obtained from any layer.
- $\Rightarrow$  We observe that the created image gets more abstract as we take embeddings from deeper layers.

Deep Dream: Suppose instead of starting with a blank image, we start with an actual image. Now we select some neuron of the network and our objective is to change the image such that this neuron fires even more. Our objective becomes  $\max L(I)$ .

Where  $h_{ij}$  is the output of the  $j^{th}$  neuron of the  $i^{th}$  layer.

We can repeat the same procedure of forward propagation and back propagation and update the image pixels  $I[i, j] = I[i, j] + \eta \nabla_{I[i, j]} L(I)$ .

• Doing this iteratively would make the image more and more like the patterns that cause the neuron to fire.

The Network has been trained to detect certain patterns which appear frequently in the training dataset. So it starts to see those patterns even when they hardly exist.

⇒ For example, if a cloud looks a little bit like a bird, the network will make it look more like a bird. This in turn will make the network recognize the bird even more strongly on the next pass, until a detailed bird appears in the image.

## Deep Art

The objective of this is to provide two images; a content image, and a style image, and then create an image such that the content of the created image is same as the content image, but it is in the style of the style image.

This is done through matching the embeddings of both the content image and style image with the created image.

For capturing the content of the image:-

$$L_{\text{content}}(\vec{p}, \vec{x}) = \sum_{ijk} (p_{ijk} - x_{ijk})^2 \quad \text{where } \vec{p} \text{ is the embeddings of the content image. } \vec{x} \text{ is the embeddings of the created image.}$$

This is done for some layers of the created image. This is done for some layers and the formula is written in the tensor notation where the final output embeddings is in the form of tensor.

For capturing the style of the image:-

It turns out that if  $V \in \mathbb{R}^{64 \times 256 \times 256}$  is the activation/output at a layer, then  $VT \in \mathbb{R}^{164 \times 64}$  captures the style of the image.

The deeper layers capture more of this style information.

The objective can be the following:-

$$L_{\text{style}} = E_e = \sum_{ij} (G_{ij} - A_{ij})^2, \text{ where } G \text{ is the Gram matrix}$$

of the style image computed at layer L; and A is similarly for the created image.

⇒ The total loss is given as:-  $L_{\text{total}}(I) = \alpha L_{\text{content}}(I) + \beta L_{\text{style}}$

### Fooling Deep CNNs

By optimizing over images, we can also fool deep CNN's to wrongly predict classes for obvious images. This is done through maximizing the log likelihood of the incorrect class and making changes in the images. The resultant images fools the network into predicting the wrong classes.

Now if instead, we gave blank image as input, then the resultant image is very noisy and un-interpretable; but, even then the network convincingly predicts it as belonging to a class. It samples little random noise & expect to get this noise from background images are very high dimensional entities. This happens because images are just vectors through a function which has its decision boundary; Even though the image is noisy to us, the CNN predicts it according to the output of the function (which is a vector function). Therefore if

## Deep Learning Part - 2

Representing Joint Distributions:- Suppose there are 3 random variables :- I takes 2 values, G takes 3 values, S takes 2 values. Now if we want to represent their joint distribution, we would require  ~~$2 \times 3 \times 2 = 12$~~   $2 \times 3 \times 2 - 1 = 11$  values, last one is 1 - (sum of 11). If we use chain rule for representing the joint representation :-  $P(S, G, I) = P(S|G, I) \times P(G|I) \times P(I)$

$P(I)$  can be represented with 1 value :-  $P(I_1)$  or  $P(I_2)$   
 $P(G|I)$  can be represented with 4 values :-  $P(G_1|I_1), P(G_2|I_1), P(G_1|I_2), P(G_2|I_2)$   
 $P(S|G, I)$  can be represented with 6 values :-  $P(S_1|G_1, I_1), P(S_2|G_1, I_1), P(S_1|G_2, I_1), P(S_2|G_2, I_1), P(S_1|G_1, I_2), P(S_2|G_2, I_2)$   
 $P(S_1|G, I)$  has 6 different combinations of  $I$  and  $G$ .  
 $P(S_1|G, I)$  and  $P(S_2|G, I)$  can be calculated/derived

In total we still take  $1+4+6=11$  values to represent the joint distribution. But if we make some assumptions on the conditional independence of S and G given I, then the number of parameters can be reduced. So  $P(S, G|I) = P(S|I) \times P(G|I)$ . Hence the joint

$P(G|S, I) = P(G|I)$  distribution can now be written as :-  $P(S, G, I) = P(S|I) \times P(G|I) \times P(I)$   
 $P(S_1|I_1)$  and  $P(S_2|I_2)$  have 2 values, 4 values and 1 value.

The number of parameters reduces to  $2+4+1=7$  values. Assumptions like this can be made based on the application, the advantages are :- 1) more natural, variables that exhibit dependence on others is made explicit. of this representation

- 2) More compact, less parameters required if conditional independence is assumed.
- 3) More modular, these values can be reused if more variables are added later, reusing values is not an option in plain representation.

Bayesian Networks :- Nodes are random variables and directed edges specify the dependence relations between the random variables.

For example :-

If,  $G$  takes 3 values and all others take 2 values, the plain representation takes = 47 value

The joint distribution can be specified as:-

$$P(D, I, G, S, L) = P(D) \cdot P(I) \cdot P(S|I) \cdot P(G|D, I) \cdot P(L|G)$$

With the compact representation takes  $1+1+2+8+3 = 15$  values.

With the graph structure along with conditional probability distributions, The graph structure along with conditional probability distributions is called a Bayesian Network.

The independence assumptions are made by us as a modelling choice and the number of parameters and data required to train the model is a consequence of our modelling choice.

Reasoning Encoded in Bayesian Network :-

i) Causal Reasoning :- we try to predict downstream effects of various factors. For example asking question like what is the chance, the student gets a good LOR ( $L'$ ) given that they are intelligent ( $I'$ ). Intelligence ( $I'$ ) is a factor and how it affects  $L$  is the question that is being asked here.  $P(L'|I')$   
 $P(L'|I') = \frac{P(L', I')}{P(I')}$ , we expect that  $P(L'|I') > P(L')$  since the chance of getting a good LOR increases given the student is intelligent.

ii) Evential Reasoning :- We reason about causes by looking at their effects. Asking questions like what is the probability of the student being intelligent given that he got a poor grade ( $G^3$ ), we expect  $P(I'|G^3) < P(I')$ , same for difficulty  $P(D'|G^3) > P(D')$ , the difficulty of course is most probably hard given the student scored poorly.  
(increases)

iii) Explaining Away :- we see how different causes of the same effect can interact. We already saw that  $P(I'|G^3)$

But what if we are given another cause of the grade, the difficulty. So probability of student being intelligent given he scored poorly and the course was difficult.  $P(I')$   
 $P(I'|G^3, D) > P(I'|G^3)$ , is our expectation since now the course was also difficult, this compensates ~~for~~ the drop in intelligence. So in a way, various causes explain the variation in other causes.

Independencies Encoded by a Bayesian Network :-

Given  $n$  random variables, we are interested in knowing if

$X_i \perp X_j$  or  $X_i \perp X_j | Z$  where  $Z \subseteq X_1, X_2, \dots, X_n / X_i, X_j$

Rule 1:- A node is not independent of its parents. (By construction)

Rule 2:- A node is not independent of its parents even when we are given the values of other variables. For example knowing the value of I does not mean that L and G become independent. L still depends on G even after knowing I.

These were some independencies between nodes and their parents. Now what about non-parents. L can still be affected by S which is a non-parent, since S affects I and I affects G which then affects L. So  $L \not\perp S$ , but given the value of G,  $L \perp S$  holds. So a node is independent of non-descendants given the parents. So  $L \perp S | G$ , moreover  $L \perp I | G$  also holds.

This is because, the effect S or I have on L is through G, but we already know G, so knowing S or I does not matter for L.

Independencies of a node and its descendants is also encoded.

Rule 3:- A node is independent to all non-descendants given the parent variables of the node.

Formally:- A Bayesian Network is a directed acyclic graph where nodes represent variables. Let  $P_{X_i}$  denote the parents of  $X_i$ .

and  $ND(X_i)$  denotes the non-descendants of  $X_i$  (does not include parents). Then for each variable  $X_i$  the graph encodes the following conditional independence assumptions:-

$$I_i(G) = \left( (X_i \perp ND(X_i)) \mid P_{X_i} \right)$$

I-Maps :- Let  $P$  be a joint distribution over  $X = X_1, X_2, \dots, X_n$ .

$I(P)$  be the set of independence assumptions that hold in  $P$ . Each element of this set is of the form  $X_i \perp X_j \mid Z, Z \subseteq X \setminus X_i$ . Let  $I(G)$  be the set of independence assumptions associated with a graph  $G$ . We say that  $G$  is an I-map for  $P$  if  $I(G) \subseteq I(P)$ . Hence all the assumptions that  $G$  states hold in  $P$ . But  $P$  can have additional independencies.

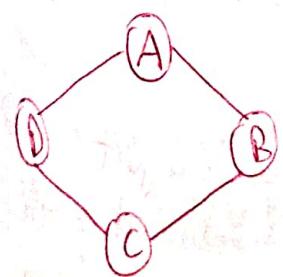
In practice, we do not know  $P$  and hence can't compute  $I(P)$ . So we just make assumptions about  $I(P)$  and then construct a  $G$  such that  $I(G) \subseteq I(P)$ .

Theorem :- Let  $G$  be a BN over a set of random variables  $X$  and let  $P$  be a joint distribution over these variables. If  $G$  is an I-Map for  $P$ , then  $P$  factorizes according to  $G$ .

Theorem :- Let  $G$  be a BN over a set of random variables  $X$  and let  $P$  be a joint distribution over these variables. If  $P$  factorizes according to  $G$ , then  $G$  is an I-Map of  $P$ .

Perfect Map :- A graph  $G$  is a Perfect Map for a distribution  $P$  if the independence relations implied by the graph are exactly the same as those implied by the distribution.

### Undirected Graphical Models :-

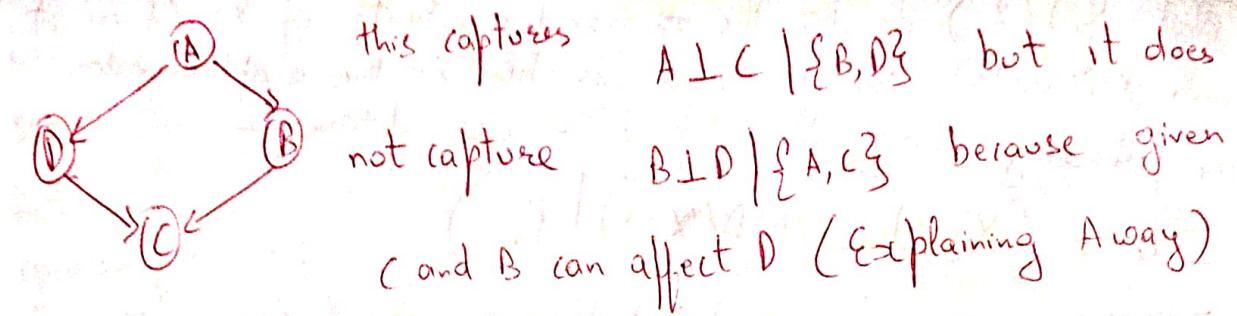


The conditional independencies in this problem are:-

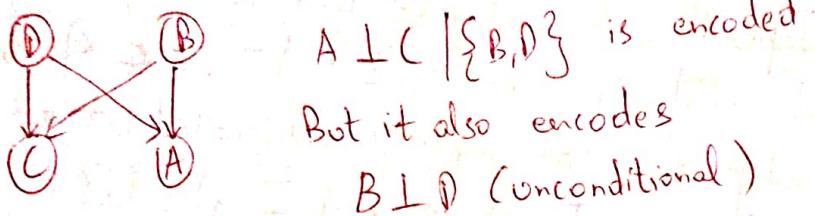
$$A \perp C \mid \{B, D\} \quad (\text{because } A \text{ and } C \text{ never interact})$$

$$B \perp D \mid \{A, C\} \quad (\text{because } B \text{ and } D \text{ never interact})$$

What if we try to model this using Bayesian Networks?



this captures  $A \perp C | \{B, D\}$  but it does not capture  $B \perp D | \{A, C\}$  because given (A and B can affect D) (Explaining Away)



$A \perp C | \{B, D\}$  is encoded

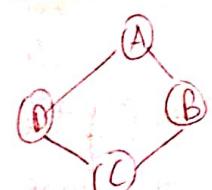
But it also encodes  
 $B \perp D$  (unconditional)

Turns out that no Bayesian Network can exactly capture independence relations completely (Perfect Map)

This is because a directed edge between two nodes implies some kind of direction in the interaction. But in our example A and B are equal partners. We instead want to capture the strength of the interaction between the nodes. These Undirected network models are called Markov Networks.

We can't use Conditional Probability Distributions as they don't make sense in the Undirected case since there is no natural conditioning.

Hence the parameters in this case is the affinity strength between connected random variables. We can have factors  $\phi_1(A, B)$ ,  $\phi_2(B, C)$ ,  $\phi_3(C, D)$ ,  $\phi_4(A, D)$  which capture the affinity between the corresponding nodes. Suppose All variables take 2 values, then:



$$\phi_1(A, B)$$

$$\begin{matrix} a & b \\ a & b \end{matrix} \quad 30$$

$$\begin{matrix} a & b \\ a & b \end{matrix} \quad 5$$

$$\begin{matrix} a & b \\ a & b \end{matrix} \quad 1$$

$$\begin{matrix} a & b \\ a & b \end{matrix} \quad 10$$

$$\phi_2(B, C)$$

$$\begin{matrix} b & c \\ b & c \end{matrix} \quad \dots$$

$$\begin{matrix} a & c \\ a & c \end{matrix} \quad \dots$$

$$\begin{matrix} a & c \\ a & c \end{matrix} \quad \dots$$

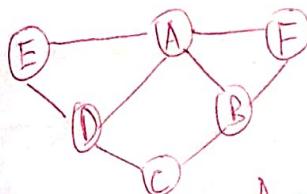
$$\begin{matrix} a & c \\ a & c \end{matrix} \quad \dots$$

Notice that here parameters are defined over the edges, while in BN's the parameters were defined over the nodes.

These tables are not probability distributions. These are just weights which can be interpreted as the relative likelihood of an event. To convert these into probability distributions we can normalize the product :-

$$P(a, b, c, d) = \frac{1}{Z} \cdot \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a)$$

where  $Z = \sum_{a, b, c, d} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a)$   
 for all values of  $a, b, c, d$ .  
 Combinations.



We could still factorize pairwise or we could factorize based on maximal cliques.  
 A maximal clique is a set of nodes in which adding any other will result in a non-clique set.

So the above network can be factorized like this:-

$$\phi_1(A, D, E) \cdot \phi_2(A, B, F) \cdot \phi_3(B, C) \cdot \phi_4(C, D)$$

A clique is a set of nodes which are pairwise connected to each other.

So a distribution factorizes over a Markov Network  $H$  if  $P$  can be expressed as  $P(X_1, \dots, X_n) = \prod_{i=1}^m \phi_i(D_i)$  where  $D_i$  is a complete sub graph (clique in  $H$ )

A distribution is a Gibbs distribution parametrized

by a set of factors  $\Phi = \{\phi_1(D_1), \dots, \phi_m(D_m)\}$  if it is defined

$$P(X_1, X_2, \dots, X_n) = \frac{1}{Z} \prod_{i=1}^m \phi_i(D_i)$$

set of all Random variables.

$\Rightarrow$  Let  $X, Y, Z$  be some distinct subsets of  $U$ , then a distribution  $P$  over these RVs would imply  $X \perp Y | Z$  if and only if we can write  $P(X) = \phi_1(X, Z) \phi_2(Y, Z)$

In our example of 4 variables,  $P(a, b, c, d) = \frac{1}{Z} \phi_1(a, b) \cdot \phi_2(b, c) \cdot \phi_3(c, d) \cdot \phi_4(d, a)$

can be written as :-  $P(a, b, c, d) = \frac{1}{Z} \phi_5(b, \{a, c\}) \cdot \phi_6(d, \{a, c\})$

Hence  $B \perp D | \{A, C\}$  holds. Similarly  $A \perp C | \{B, D\}$  also holds.

For a given Markov Network  $H$  we define Markov Blanket of a RV  $X$  to be the neighbors of  $X$  in  $H$ .

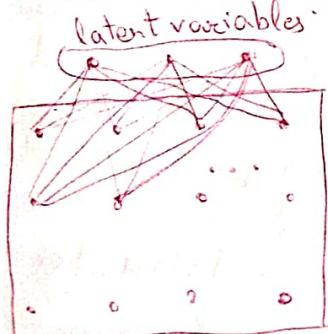
Analogous to the case of Bayesian Networks we can define the local independencies associated with  $H$  to be :-

$$X \perp (U - \{x\} - MB_H(X)) \mid MB_H(X)$$

$X$  is independent of its non-neighbors given the neighbors of  $X$  in  $H$ . The neighborhood is defined by the Markov Blanket. For example in an image, the Markov Blanket can be the 8 pixels that surround it.

Latent Variables:- Now in the case of modelling an image as a Markov Network, each pixel can be treated as a random variable with blanket being the 8 surrounding pixels. So the joint distribution of images can be  $P(V_1, V_2, \dots, V_{1024})$  where  $V_i$  is a pixel and can be factorized accordingly.

Instead of just modelling the pixels as random variables, we can introduce latent variables such that the value of pixels depend on the ~~value~~ latent variables. Consider the task of generating a scenery. The latent variables can be:- is the scene sunny or not, does the scene have mountains or beaches, clouds or not etc. Based on these variables, the actual pixels can be generated. The Markov Network looks like this:-



pixels are dependent on the latent variables.  
The interactions between the pixels are captured through the latent variables.

This modelling choice helps in many tasks like abstraction. Suppose, we are able to learn the joint distribution  $P(V, H)$ . Using this distribution we can find

$$P(H|V) = \frac{P(V, H)}{\sum P(V, H)} \rightarrow \text{marginalize over } V. \text{ So given an image,}$$

we can find the most likely latent configuration of  $H$  that generated this image. The vector  $h$  calculated can be thought of as an abstract representation of the whole image. Similar images which are pixel wise still very different in the original space are expected to have very similar latent configuration.

→ In practice we don't have latent configurations given to us. This is like a Representation Learning Task and each element in the latent vector representation may not even be interpretable. They just help us to learn a good abstraction of the data.

Another task these representations can help is ~~is~~ generation. Using the joint distribution, we can find:  $P(V|H) = \frac{P(V, H)}{\sum P(V, H)}$

So, given a latent configuration, generate the image which most likely ~~has~~ <sup>corresponds to this</sup> configuration.

Let the Vector  $V$  of all visible variables be a binary vector,  $V \in \{0, 1\}^m$  and the vector  $H$  also be a binary vector

$$H \in \{0, 1\}^n$$

Notice that edges are only incident between each pair of hidden-visible variables. We do not have edges between (hidden-hidden) and (visible-visible) variables. Hence the cliques in the graph only includes each pair of (hidden-visible) pair of variables. So the joint distribution can be factorized as:  $P(V, H) = \frac{1}{Z} \prod_{i=1}^m \prod_{j=1}^n \phi_{ij}(v_i, h_j)$

We can add additional factors such that

$$P(V, H) = \frac{1}{Z} \prod_{i=1}^m \prod_{j=1}^n \phi_{ij}(v_i, h_j) \prod_{i=1}^m \psi_i(v_i) \prod_{j=1}^n \epsilon_j(h_j) \quad \text{this is}$$

allowed as long as  $Z$  correctly normalizes to obtain a prob distribution.  $Z$  is calculated as follows:

$$Z = \sum_V \sum_H \prod_{i=1}^m \prod_{j=1}^n \phi_{ij}(v_i, h_j) \prod_{i=1}^m \psi_i(v_i) \prod_{j=1}^n \epsilon_j(h_j)$$

all  $2^m$  possible configurations  $\rightarrow$  all  $2^n$  possible configurations. Hence this  $Z$  term includes  $2^{m+n}$  terms in total.

Each  $\phi_{ij}$  is a table which gives the values for each possible (4) configuration of  $v_i$  and  $h_j$ . Similarly  $\psi_i$  and  $\epsilon_j$  are also tables giving values for each possible (2) configuration of  $v_i$  and  $h_j$  respectively.

Restricted Boltzmann Machines :- TD learn  $\phi_{ij}$ 's, ~~and~~  $\psi_i$ 's and  $\epsilon_j$ 's, we introduce parameters. RBMs chooses the

following parametric forms :-  $\phi_{ij}(v_i, h_j) = e^{w_{ij} v_i h_j}$

With this parametric form,  $\psi_i(v_i) = e^{b_i v_i}$

The distribution can be

expressed as follows :-

$$P(V, H) = \frac{1}{Z} \prod_{i=1}^m \prod_{j=1}^n e^{w_{ij} v_i h_j} \prod_{i=1}^m b_i v_i \prod_{j=1}^n e^{c_j h_j}$$
$$= \frac{1}{Z} e^{\left( \sum_{i=1}^m \sum_{j=1}^n w_{ij} v_i h_j + \sum_{i=1}^m b_i v_i + \sum_{j=1}^n c_j h_j \right)}$$

$$P(V, H) = \frac{1}{Z} e^{-E(V, H)}$$

where  $E(V, H) = -\sum_{i=1}^m \sum_{j=1}^n w_{ij} v_i h_j - \sum_{i=1}^m b_i v_i$   
 $- \sum_{j=1}^n c_j h_j$

$E(V, H)$  is called the energy

High energy  $\rightarrow$  low probability, and vice versa.

This is called RBM because, we do not allow edges between the visible vertices (restricted), and the form of joint distribution is similar to Boltzmann/Gibbs distribution.

RBM's as stochastic Neural Networks :- Let us try to model

$P(H|V)$  and  $P(V|H)$  given the form of joint distribution.

Let  $V_{-l}$  be the state of all visible units except the  $l^{th}$  unit.  
 $\alpha_l(H) = -\sum_{j=1}^n w_{ej} h_j - b_e$  part of  $l^{th}$  visible unit in the summation of  $E(V, H)$ .

$$\beta(V_{-l}, H) = -\sum_{i=1}^m \sum_{j=1, j \neq l}^n w_{ij} v_i h_j - \sum_{i=1}^m b_i v_i - \sum_{j=1}^n c_j h_j$$

$$\text{So } E(V, H) = V_l \alpha_l(H) + \beta(V_{-l}, H)$$

We can hence calculate  $P(V_l = 1 | H)$ . Notice that since  $V_{-l}$  is independent of  $V_l$  given  $H$ .  $P(V_l = 1 | H) = P(V_l = 1 | V_{-l}, H)$  doesn't matter what configuration  $V_{-l}$  is.

Representing it in this way allows for using the joint distribution :-  $P(V_e=1|H) = P(V_e=1|V_{-e}, H) = \frac{P(V_e, V_{-e}, H)}{P(V_{-e}, H)}$

$$- E(V_{\ell=1}, V_{-\ell}, H) \quad p(V_{-\ell}, H)$$

$$\Rightarrow \frac{e}{-E(V_L=0, V_{-L}, H)} + \frac{-E(V_L=1, V_{-L}, H)}{e} = -\alpha e^{-H} - \beta(V_{-L}, H)$$

$$e^{-\alpha_l^{(n)} \beta(V_{-l}, H)} + e^{-\beta(V_{-l}, H)}$$

$$\Rightarrow \frac{e^{-\beta(V_{-l}, H)} - \alpha_l(H)}{e^{-\beta(V_{-l}, H)} + e^{-\alpha_l(H)}} = \frac{e^{-\alpha_l(H)}}{e^{-\alpha_l(H)} + 1} = \frac{1}{e^{\alpha_l(H)} + 1}$$

$$\Rightarrow P(v_e=1 | H) = \sigma(-\alpha_e(H)) = \sigma\left(\sum_{j=1}^n w_{ej} h_j + b_e\right)$$

Similarly it can be shown that  $P(h_p=1|V) = \sigma\left(\sum_{i=1}^n w_{ip} V_i + b_p\right)$

<sup>extra</sup>  
Hence, the factors added to the form of  $p(H; V)$  serve as biases in the formulation of RBMs. So RBM's can be interpreted as stochastic neural network where the nodes and edges correspond to neurons and synaptic connections, respectively. The conditional probability of a single variable being 1 can be interpreted as the firing of a neuron with sigmoid activation.

Having  $P(y_0=1|H)$ , the probability  $P(V|H)$  is straightforward.

$$P(V_i = v_i | H) = \prod_{i=1}^m P(V_i = v_i | H) \quad \begin{matrix} \text{since } V_i \text{ are components} \\ \text{independent independent} \end{matrix}$$

Since our ~~dedicated~~<sup>task</sup>, ~~there~~ is

Since our objective task is to find the solution

Unsupervised, our objective function should only depend on  $V$ 's  $\rightarrow$  the images. Eventually we want to

Should only depend on  $V_S$ .  
 maximize the likelihood of our data. So the task is now to  
 find parameters  $w_{ij}, b_i, \alpha_j$  that maximize the probability of

the images. So we want to maximize  $P(V|O)$ , the set of parameters  $w_{ij}, b_i, w_j$ .

So maximize  $\prod_{i=1}^m P(V_i=v_i|\theta)$ . Lets consider the loss for one single example.  $P(V_i=v_i)$  is same as  $P(V_i=v_i|\theta)$

$$\text{So } P(V) = \prod_{i=1}^m P(V_i=v_i)$$

$$\text{Also we know that } P(V, H) = \frac{1}{Z} e^{-E(V, H)} \text{ So } P(V) \text{ can also}$$

$$\text{be calculated } \Rightarrow \frac{1}{Z} \sum_H e^{-E(V, H)}$$

→ marginalized over  $H$   
all possible  $2^n$  configurations

$$\text{So } P(V) = \frac{1}{Z} \sum_H e^{-E(V, H)} = \frac{\sum_H e^{-E(V, H)}}{\sum_{V, H} e^{-E(V, H)}} \text{, marginalized over } 2^{n+m} \text{ configurations of } V \text{ and } H.$$

$$\text{So } \ln P(V) = \ln \sum_H e^{-E(V, H)} - \ln \sum_{V, H} e^{-E(V, H)}$$

We maximize  
this instead

$$\frac{\partial \ln P(V)}{\partial \theta \text{ any param}} = -\frac{1}{\sum_H e^{-E(V, H)}} \sum_H e^{-E(V, H)} \cdot \frac{\partial E(V, H)}{\partial \theta} + \frac{1}{\sum_{V, H} e^{-E(V, H)}} \sum_{V, H} e^{-E(V, H)} \cdot \frac{\partial E(H, V)}{\partial \theta}$$

$$= -\sum_H e^{-E(V, H)} \left\{ \frac{\partial E(V, H)}{\partial \theta} \right\} + \sum_{V, H} e^{-E(H, V)} \left\{ \frac{\partial E(H, V)}{\partial \theta} \right\}$$

after dividing both  
numerators and denominators  
by  $Z$ , observe that

$$\frac{P(V, H)}{P(V)} = P(H|V)$$

$$\text{Hence } \boxed{\frac{\partial \ln P(V)}{\partial \theta} = -\sum_H P(H|V) \cdot \frac{\partial E(V, H)}{\partial \theta} + \sum_{H, V} P(V, H) \cdot \frac{\partial E(H, V)}{\partial \theta}}$$

$$\text{if } \theta = w_{ij} \text{ then } \frac{\partial L(\theta)}{\partial w_{ij}} = + \sum_H p(H|V) h_j v_i - \sum_{H,V} p(V,H) h_j v_i \\ = + \underset{p(H|V)}{\mathbb{E}[h_j v_i]} - \underset{p(H,V)}{\mathbb{E}[h_j v_i]}$$

Both these expectations have an exponential number of terms. And notice that this is gradient wrt 1 param of a single observation. Hence the learning is intractable in its current form. So we need sampling methods to estimate these expectations. But sampling is not straightforward as only a few samples are likely from the huge space of samples. So Sampling has to be on the basis of the joint distribution but we don't have the distribution itself.

Solution is to draw samples from an easier distribution as long as sampling from the other distribution is same as if they were drawn from  $P(\text{original})$ . We do this in Gibbs Sampling.

Goal 1: Given a random variable  $X \in \mathbb{R}^n$  we are interested in drawing samples from the joint distribution  $P(X)$

Goal 2: Given a function  $f(X)$  defined over the random variable  $X$ , we are interested in computing the expectation  $E_{P(X)}[f(X)]$

Suppose we have a chain of random variables

$X_1, X_2, \dots, X_n$  each  $X_i \in \mathbb{R}^n$ , basically each random variable  $X_i$  is in itself a vector.  $i$  corresponds to a time step. For our discussion, let  $X_i \in \{0, 1\}^n$ . One way of looking is to see that  $X_i$  transitions to  $X_{i+1}$  when the state changes from  $i$  to  $i+1$ .  $X_i$  can take one of  $2^n$  values in each state.

One might be interested in finding the most likely value  $X_i$  will take given all the values of previous states  $X_{i-1}, X_{i-2}, \dots, X_1$ .

So we want :-  $P(X_i = x_i | X_1 = x_1, X_2 = x_2, \dots, X_{i-1} = x_{i-1})$

But let the chain exhibit the Markov Property:-

$$P(X_i = x_i | X_1 = x_1, X_2 = x_2, \dots, X_{i-1} = x_{i-1}) = P(X_i = x_i | X_{i-1} = x_{i-1})$$

given the previous state  $X_{i-1}, X_i$  is independent of all preceding states.  $\textcircled{x}_1 \rightarrow \textcircled{x}_2 \rightarrow \dots \rightarrow \textcircled{x}_k$   $[X_i \perp\!\!\! \perp X_1^{i-2} | X_{i-1}]$ .

This chain of random variables is called a Markov Chain.  
We are dealing with a discrete time and discrete space Markov Chain.  
 $i$  is discrete  $X_i$  takes discrete values.

Let  $2^n = l$ . Now to specify  $P(X_i = x_i | X_{i-1} = x_{i-1})$  we need to specify  $l^2$  values, for each configuration of  $X_{i-1}$  and  $X_i$ . We can represent this as a matrix  $T \in l \times l$  where  $T_{ab}$  gives the probability of transitioning from  $X_{i-1} = a$  to  $X_i = b$ , where 'a' and 'b' encode ~~two~~ of the possible  $2^n$  configurations respectively.  $T$  is called the transition matrix.  $T$  may be different for each transition, but we make the assumption that our Markov Chain is time homogeneous. So for all transitions,  $T$  remains the same.  $P(X_i = b | X_{i-1} = a) = T_{ab}$   $\forall a, b, i$

Let the distribution be specified by  $u^k$  at  $k^{\text{th}}$  time step.  $u^k$  is a  $l$  dimensional vector.  $u_a^k = P(X_k = a)$ . ~~Each~~ The  $a^{\text{th}}$  element of  $u^k$  gives the probability of  $X_k$  being 'a'. Again,  $X_k$  is an  $n$ -dimensional vector but 'a' is an encoding of one of such possible vector.

Let the starting distribution be  $u^0$ . How to calculate  $P(X_i = b)$  from  $u^0$ ?  $P(X_i = b) = \sum_a P(X_0 = a, X_i = b)$  marginalize over  $X_0$ .

Probability of  $X_i$   
taking 'b' irrespective  
of the starting point

$$= \sum_a P(X_0 = a) \cdot P(X_i = b | X_0 = a)$$

$$= \sum_a u_a^0 T_{ab} = u_b^i \text{ since } P(X_i = b) = u_b^i$$

So  $u^i$  can be written as  $u^i = u^0 T = [u_1^0, u_2^0, \dots, u_l^0] \begin{bmatrix} T_{11} \\ T_{21} \\ \vdots \\ T_{l1} \end{bmatrix}$

Similarly  $u^2 = u^0 T^2$

So in general  $u^k = u^0 T^k$  this gives us a way to get the distribution at time  $k$ . To get the random variables at  $k$ , we can perform sampling from  $u^k$ .

This is still computationally expensive as the dimensions of the vectors involved are exponential in terms of the number of variables.

Also notice that if at some time step  $t$ ,  $\pi^t$  reaches a distribution  $\pi$  such that  $\pi^T = \pi^t = \pi$ . Then for all subsequent time steps  $\pi^j = \pi$  for  $j \geq t$ .

$\pi$  is called the stationary distribution of the Markov Chain

$X_t, X_{t+1}, X_{t+2}, \dots$  will all follow the same distribution  $\pi$ .

So if we run a Markov Chain for a large number of time steps then after a point we start getting samples  $x_t, x_{t+1}, x_{t+2}$ ,

which are essentially being drawn from the stationary distribution.

But as mentioned earlier, drawing from this distribution is computationally expensive and still intractable.

But suppose if we are given ~~a~~ a Markov Chain such that it is easy to draw samples from it and the stationary distribution of this chain is  $P(X)$ , Then we could empirically

estimate  $E_{P(X)}[f(x)]$  as  $\frac{1}{n} \sum_{i=t}^{t+n} f(x_i)$

$\rightarrow$  the time step after we get the stationary distribution.

Theorem:-

If  $X_0, X_1, \dots, X_t$  is an irreducible time homogeneous discrete Markov Chain with stationary distribution  $\pi$ , then

$$\frac{1}{t} \sum_{i=1}^t f(x_i) \xrightarrow[\text{almost surely}]{} E_{\pi}[f(x)] \text{ where, } X \in \mathcal{X} \rightarrow \{0, 1\}^n$$

for any function  $f: \mathcal{X} \rightarrow \mathbb{R}$   $X \sim \pi$  in our case

$T_{10}$  } further the Markov Chain is aperiodic then  
 $T_{01}$  }  
 $T_{11}$  }

$$P(X_t = x_t | X_0 = x_0) \rightarrow \pi(x) \text{ as } t \rightarrow \infty \quad \forall x, x_0 \in \mathcal{X}$$

Basically means that doesn't matter the starting state, after large number of steps, the samples calculated are drawn from  $\pi(x)$ .

Now our task is to find such a Markov Chain whose  $\pi(x)$  is same as  $p(x)$ , moreover it is irreducible and aperiodic.

For ease of notation, instead of  $X = V_1, V_2, \dots, V_m, H_1, H_2, \dots, H_n$  we will use  $X = X_1, X_2, \dots, X_{n+m}$  (Basically renaming them all to  $X$ ).

We will now set up a Markov Chains for RBMs.

Here  $X_i$  does not mean random vector at time step  $i$ , here it means the  $i^{\text{th}}$  element of the random vector, basically a scalar 0 or 1.

We set up the chain in this way:- i) We start at a random  $X$  where each element is set to 0 or 1 uniformly at random.

ii) Sample a value  $i \in \{1, \dots, n+m\}$  (index) using a distribution,  $q(i)$ , can be uniform.

iii) Fix the value of all variables except  $X_i$ . Sample a new value

for  $X_i$  using the following conditional distribution

$$P(X_i=y_i | X_{-i}=x_{-i}) \quad \text{conditioned on remaining variables.}$$

So a transition from state  $x$  to state  $y$  is only possible if  $x$  and  $y$  differs at most by 1 variable out of  $n+m$  variables. All other transitions are not possible in this chain. Hence the transition matrix is extremely sparse and can be explained like this:-

$$T_{xy} = \begin{cases} q(i) P(X_i=y_i | X_{-i}=x_{-i}), & \text{if } \exists i \in \{1 \dots n+m\} \\ 0, & \text{otherwise} \end{cases} \quad \text{such that if } v \neq i \quad x_v = y_v$$

Observe that in the case of RBMs,  $P(X_i=y_i | X_{-i}=x_{-i})$  is easily calculated. If  $i \leq m$  (visible variable is selected)

$$P(V_i=1 | V_{-i}, H) = P(V_i=1 | H) = \sigma \left( \sum_{j=1}^n w_{ij} h_j + b_i \right)$$

and if  $i > m$ , (hidden variable is selected)

$$P(H_j=1 | V, H_{-j}) = P(H_j=1 | V) = \sigma \left( \sum_{i=1}^m w_{ij} V_i + c_j \right)$$

Hence the sampling procedure at each step is as follows:-

1.  $\Rightarrow$  Sample  $i$  from  $q(i)$

2.  $\Rightarrow$  Sample  $X_i$  from a Bernoulli Distribution, keeping all variables the sample same. Both these computations are easy. Hence it is easy to sample from the chain.

3. Theorem (Detailed Balanced Condition) :- To show that a distribution  $\pi$  is a stationary distribution for a Markov Chain described by the transition probabilities  $T_{xy}$ , it is sufficient to show that  $\forall x, y$  the following condition holds:-

$$\boxed{\pi(x) T_{xy} = \pi(y) T_{yx}}. \text{ So in our case if we want}$$

to prove that  $P(X)$  is actually the stationary distribution of this chain then it suffices to show that  $P(x) T_{xy} = P(y) T_{yx}$

To prove this, we look at 3 cases, when  $x$  and  $y$  differ in more than variables, 1 variable, are same - i) differ in more than 2

variables:- In this case  $T_{xy} = 0$  And  $T_{yx} = 0$ , hence

$$P(x) T_{xy} = P(y) T_{yx} = 0 \text{ holds true}$$

ii) when  $x$  any  $y$  are same so  $P(x) = P(y)$  and  $T_{xy} = T_{yx}$

$$\text{So } P(x) T_{xy} = P(x) T_{xx} = P(y) T_{yx}, \text{ holds true.}$$

iii) when  $x$  and  $y$  differ in 1 variable then, let that variable be the  $i^{th}$  variable.

$$\text{So, } T_{xy} = q(i) P(x_i = y_i | X_{-i} = x_{-i})$$

$$P(x) T_{xy} = P(x_i, x_{-i}) \cdot q(i) \cdot \frac{P(y_i, x_{-i})}{P(x_{-i})}$$

$$= P(y_i, x_{-i}) \cdot q(i) \cdot \frac{P(x_i, x_{-i})}{P(x_{-i})}$$

Since other variables  
are same

$$= P(y) \cdot q(i) \cdot P(x_i = x_i | X_{-i} = x_{-i})$$

$$P(x) T_{xy} = P(y) \cdot T_{yx} \quad [ \because L \rightarrow = P(x_i = x_i | X_{-i} = y_{-i}) ]$$

Hence the detailed balance condition holds and  $P(X)$  is indeed the stationary Distribution of this chain.

Now we need to prove that this chain is irreducible and aperiodic.

A markov chain is irreducible if one can get from any state in the sample space to any other state in a finite number of transitions.

For  $\forall y \exists k > 0$  with  $P(X^{(k)} = y | X^{(0)} = x) > 0$ .

Our Markov chain is irreducible as from any starting state, any other state ~~is~~ only differs in finite number of variables, and those variables can be flipped in one by one with a non zero probability.

A chain is called aperiodic if ~~forall~~ in the sample space, the greatest common divisor of the set :-

$\{k \mid P(X^{(k)} = x | X^{(0)} = x) > 0\}$  is 1. This set is the set

of all those time steps such that the initial value is obtained again. Having a GCF of 1 essentially means that these numbers are arbitrary and ~~not~~ are not a multiple of any number which would otherwise make the chain periodic. This has to satisfy for all  $x$  in the sample space. Our Markov chain is aperiodic since the initial value can be retained in each time step with a non-zero prob, and gcd of all natural numbers is 1.

We return to gradient computation of RBMs.

$$\frac{\partial L(\theta)}{\partial w_{ij}} = \frac{E[h_j v_i]}{P(H|V)} - \frac{E[v_i h_j]}{P(V|H)} = \sum_H p(h|v) h_j v_i - \sum_{v,h} p(v,h) h_j v_i \\ = \sum_H p(h|v) h_j v_i - \sum_v \sum_H p(H|V) h_j v_i$$

$$\sum_H p(h|v) h_j v_i = \sum_{H_j} \sum_{H-j} p(h_j|v) p(H-j|v) h_j v_i \\ = \sum_{H_j} p(h_j|v) h_j v_i \left( \sum_{H-j} p(H-j|v) \right) \xrightarrow{\text{this just sum of } n-1 \text{ values here. Hence a distribution. Hence equal to 1.}} \\ = p(h_j=1|v) v_i + p(h_j=0|v) v_i \times 0$$

$$= p(h_j=1|v) v_i = \sigma \left( \sum_{i=1}^m w_{ij} v_i + (j) \right) v_i$$

$$\text{Hence } \frac{\partial L(\theta)}{\partial w_{ij}} = \sigma \left( \sum_{i=1}^m w_{ij} v_i + (j) \right) v_i - \sum_v p(v) \sigma \left( \sum_{i=1}^m w_{ij} v_i + (j) \right) v_i$$

We can write this expression as a general form for the whole matrix  $W$ .

$$\frac{\partial L(\theta)}{\partial W} = V \cdot \sigma(V^T W + C^T) - \sum_v p(v) V \cdot \sigma(V^T W + C^T)$$

where  $V$  is the vector  $[v_1, v_2, \dots, v_m]$  and  $W$  is an  $m \times n$  matrix.  $\sigma(V^T W + C^T)$  is an element wise sigmoid of the vector inside.

$$V \cdot \sigma(V^T W + C^T) = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix}_{m \times 1} \cdot \sigma \left( [v_1, v_2, \dots, v_m]^T \begin{bmatrix} w_{11} & w_{12} & \dots & w_{1n} \\ w_{21} & w_{22} & \dots & w_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m1} & w_{m2} & \dots & w_{mn} \end{bmatrix}_{m \times n} \right)$$

Element wise Softmax operation

$$+ [c_1, c_2, \dots, c_n]$$

$$= \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_m \end{bmatrix} \cdot \sigma \left( \left[ \sum_{i=1}^m w_{i1} v_i + c_1, \sum_{i=1}^m w_{i2} v_i + c_2, \dots, \sum_{i=1}^m w_{in} v_i + c_n \right] \right)$$

$$= \begin{bmatrix} v_1 \cdot \sigma[z_1, z_2, \dots, z_n] \\ v_2 \cdot \sigma[z_1, z_2, \dots, z_n] \\ \vdots \\ v_m \cdot \sigma[z_1, z_2, \dots, z_n] \end{bmatrix}_{m \times n}$$

Now let us derive gradients with respect to other variables:-

$$\begin{aligned} \frac{\partial L(\theta)}{\partial b_i} &= E[v_i] - E[v_i] = \frac{1}{P(H|V)} \sum_h p(h|v) v_i - \sum_{v,h} p(h,v) v_i \\ &= v_i \sum_h p(h|v) - \sum_v v_i p(v) \sum_h p(h|v) \end{aligned}$$

$$\frac{\partial L(\theta)}{\partial b_i} = v_i - \sum_v v_i p(v), \quad \frac{\partial L(\theta)}{\partial b} = V - \sum_v p(v) V$$

Similarly  $\frac{\partial L(\theta)}{\partial b_j} = E[h_j] - E[h_j]$

$$\frac{\partial L(\theta)}{\partial b_j} = P(h_j=1|v) - \sum_v p(v) P(h_j=1|v)$$

$$\frac{\partial L(\theta)}{\partial b_j} = \sigma \left( \sum_{i=1}^m w_{ij} v_i + c_j \right) - \sum_v p(v) \sigma \left( \sum_{i=1}^m w_{ij} v_i + c_j \right)$$

$$\nabla_c L(\theta) = \sigma(W^T v + c) - E_{P(v)} [\sigma(W^T v + c)]$$

Notice that all the gradient computations still involve this expectation term, which is intractable. So we will set up a markov chain and sample to estimate these expectations.

RBM training with Block Gibbs Sampling:

for all  $v_d \in D$ :  $\xrightarrow{\text{a sample}}$  set of samples.

Randomly initialize  $v^{(0)}$

for  $t=0, 1, \dots, k, k+1, \dots, K+r$ :

for  $j=1, 2, \dots, n$

Sample  $h_j^{(t)} \sim P(h_j | v^{(t)})$

for  $i=1, 2, \dots, m$

Sample  $v_i^{(t+1)} \sim P(v_i | h^{(t)})$

$$\{P(h_j=1 | v)_m =$$

$$\left( \sum_{i=1}^n w_{ij} v_i + b_j \right)$$

Notice that instead of going one by one from a sample to another, we

make steps in blocks (first all hidden, then all visible) because we have either too many time steps or we have to run the chain for large iterations. Hence each time step  $t$  actually consists of multiple time steps.

~~REWRITE~~  $h^{(t)}(W)$

$$W = W + \eta \left[ V_d \cdot \sigma(V_d^T W + c) - \frac{1}{r} \sum_{p=k+1}^{K+r} V^{(p)} \cdot \sigma(V^{(p)T} W + c) \right]$$

notice that these terms are based on the actual sample  $V_d$

these terms are based on the samples of the Markov Chain.

$$b = b + \eta \left[ V_d - \frac{1}{r} \sum_{p=k+1}^{K+r} V^{(p)} \right]$$

$$c = c + \eta \left[ \sigma(W^T V_d + c) - \frac{1}{r} \sum_{p=k+1}^{K+r} \sigma(W^T V^{(p)} + c) \right]$$

In practice, Gibbs Sampling can also be very inefficient because for every step of stochastic gradient descent, we need to run the Markov Chain for many steps.

In practice we use  $\kappa$ -contrastive Divergence to train RBMs.

Contrastive Divergence uses the following idea:-

→ Instead of starting the Markov Chain at a random point, start from  $v_d$  where  $v_d$  is the current training instance.

→ Run Gibbs Sampling for  $k$  steps, and denote the sample at the  $k^{\text{th}}$  step by  $\tilde{v}$ . Replace the expectation by a point estimate

$$E_{p(v, h)}[v_i h_j] \approx \cancel{\sigma(v_i^T W + c)} \sigma(w_j^T \tilde{v} + c_j) \cdot \tilde{v}_i$$

$\downarrow$   
 $j^{\text{th}}$  column of  $W$

Hence the gradients are as follows:-

$$W = W + \eta [v_d \cdot \sigma(v_d^T W + c^T) - \tilde{v} \cdot \sigma(\tilde{v}^T W + c^T)]$$

$$b = b + \eta [v_d - \tilde{v}]$$

$$c = c + \eta [\sigma(W^T v_d + c) - \sigma(W^T \tilde{v} + c)]$$

The higher the value of  $k$ , the less biased the estimate of the gradient will be.

Variational Autoencoders :-

before VAE's some info about autoencoders. An autoencoder contains an encoder which takes the input  $X$  and maps it to a hidden representation, then the decoder takes this hidden representation and tries to reconstruct the input from it as  $\hat{X}$ . The loss between  $\hat{X}$  and  $X$  is used to train the params of the network (Unsupervised task). Can we use Autoencoders for generative purposes? Well given a hidden representation, we can produce an  $\hat{X}$ . But it is a deterministic procedure so variability will be less. Moreover, out of all the possible hidden representations, only a few are likely to be produced by our data. So deciding ~~the~~ a good representation (hidden) for generation is still a question. In case of RBM's we have an approximate joint dist so we can find the likely hidden representations.

Variational auto encoders ~~have~~ have similar/same structure as autoencoders, but they learn a distribution over the hidden variables.

visible variables will be referred as  $X$ , and hidden variables will be referred as  $Z$ .

The architecture and goals are as follows :-

Reconstruction:  $\hat{X}$

Decoder  $P_\phi(X|z)$

$\uparrow$

$\uparrow$

Encoder  $Q_\theta(z|x)$

$\uparrow$

Data:  $x$

$\theta$ : the parameters of the encoder network.

$\phi$ : the parameters of the decoder network.

Notice that both the encoder and decoder learn distributions rather than the actual hidden/visible representations.

In the encoder part, we assume that the latent variables come from a standard normal distribution  $N(0, I)$  and the job of the encoder is to predict the parameters of this distribution. The input

of the encoder is  $x$ , and it has to give the parameters  $(\mu, \Sigma)$  of the distribution  $Q_\theta(z|x) \sim N(\mu, \Sigma)$ , and the reference distribution is  $N(0, I)$ .

The job of the decoder is to predict a probability over  $X$ :  $P(X|z)$ . We assume that  $P(X|z)$  is a Gaussian Distribution with unit variance. The decoder predicts the mean of the distribution,  $P_\phi(X|z)$ .

$\uparrow$

Decoder  $P_\phi(X|z)$

$\uparrow$

Sample  $z$

$\uparrow$

$\mu$

$\uparrow$

Encoder  $Q_\theta(z|x)$

$\uparrow$

$x$

Goal 1: - Learn a Distribution over the latent variables  $Q(z|x)$ .

Goal 2: - Learn a distribution over the visible variables  $P(x|z)$ .

Notice that both the encoder and decoder learn distributions rather than the actual hidden/visible representations.

In the encoder part, we assume that the latent variables come

from a standard normal distribution  $N(0, I)$  and the job of the encoder is to predict the parameters of this distribution. The input

of the encoder is  $x$ , and it has to give the parameters  $(\mu, \Sigma)$  of the distribution  $Q_\theta(z|x) \sim N(\mu, \Sigma)$ , and the reference distribution is  $N(0, I)$ .

The job of the decoder is to predict a probability over  $X$ :  $P(X|z)$ . We assume that  $P(X|z)$  is a Gaussian Distribution with unit variance. The decoder predicts the mean of the distribution,  $P_\phi(X|z)$ .

The loss function is as follows:-  
we want to maximize the likelihood of  $P_\phi$  the probability of the data. Basically

$$P(x_i) = \int P(z) P(x_i|z) dz$$

$$= E[\log(P_\phi(x_i|z))]$$

$$\text{Z} \sim Q(z|x)$$

$$P(z) \text{ only.}$$

log is taken for numerical stability.

But we want to also ensure that  $Q_\phi(z|x_i)$  does not diverge much from  $N(0, I)$ . So a KL divergence term between these two distributions is also taken. So the total loss :-

$$L(\phi) = -\mathbb{E}_{z \sim Q_\phi(z|x_i)} [\log(p(x_i|z))] + \text{KL}(Q_\phi(z|x_i) || N(0, I))$$

We need to add this divergence term so that the network does not cheat. Otherwise it might just learn a unique representation for each datapoint and decode it ~~as~~ as a mapping without learning anything. So we predict a distribution  $Q_\phi(z|x_i)$  such that a sample from this distribution should be able to reconstruct the original datapoint with high probability. We assume that  $p(z) \sim N(0, I)$  because ~~any~~ any distribution in d dimension ~~may~~ be generated with a set of d normally distributed variables by using a sufficiently complex function. So the decoder learns this function to arrive at  $p_\phi(x_i|z)$  using only normally distributed variables.

It can be shown that  $\text{KL}[N(u(x), \Sigma(x)) || N(0, I)]$

$$= \frac{1}{2} (\text{tr}(\Sigma(x)) + u(x)^T [u(x) - R])$$

where  $\Sigma(x)$  and  $u(x)$

are the  $u$  and  $\Sigma$  calculated by the decoder for  $\mathbb{E}_{z \sim Q_\phi(z|x_i)} [\log p(x_i|z)]$ , it is expressed like this since it is a function of  $x$ .  $R$  is the dimensionality of the latent variables. This is easily computable. But  $\mathbb{E}_{z \sim Q_\phi(z|x_i)} [\log p(x_i|z)]$  is still a

bottleneck since it involves a multivariate integral. But in VAE's we approximate this expectation using a single  $z$  sampled from  $N(u(x), \Sigma(x))$ . Hence  $\log p(x_i|z) = (-\frac{1}{2} \|x_i - \mathbb{E}_{z \sim Q_\phi(z|x_i)}[z]\|^2)$

$u(z)$  is the mean of the distribution  $\mathbb{E}_{z \sim Q_\phi(z|x_i)}[z]$  calculated by the decoder. It is of this form as  $u(z)$  depends on  $z$ , moreover the expression can be derived from the PDF of Normal distribution

Hence the effective objective function is :-

$$\text{minimize}_{\phi} \sum_{n=1}^N \left[ \frac{1}{2} [\text{tr}(\Sigma(x)) + u(x)^T [u(x) - R]] - \log \det(\Sigma(x)) + \|x - u(z)\|^2 \right]$$

Observe that  $u(x)$  and  $\Sigma(x)$  depend on  $\phi$ , while  $u(z)$  depends on  $\phi$ .

$\hookrightarrow$  we assume  $p(x|z)$  is normal with unit variance.

But since the network involves sampling  $z$ . It is not end-to-end differentiable. To make it differentiable we sample  $a$  in the input step itself. So we ~~can~~ calculate  $z$  like this --

Sample  $\epsilon \sim N(0, I)$  in the input stage. Then

$$z = u(x) + A(x)\epsilon \quad , \text{where } A(x) \text{ is such that} \\ \cancel{A(x)} \Sigma(x) = A(x)A(x)^T$$

This is because  $z$  will now follow  $N(u(x), AIA^T)$   
 $= N(u(x), \Sigma(x))$

Now the parameters  $\theta$  and  $\phi$  can be learnt using gradient descent and backpropagation.

In abstraction, we calculate  $z$  in the above way.

In generation, we sample  $z \sim N(0, I)$  ~~from~~ and just feed it to the decoder, since model was trained to converge  $\theta(z|x) \sim N(0, I)$

Autoregressive Models - these models do not contain any latent variables. They also learn a joint distribution over  $X$ . For us,  $X \in \{0, 1\}^n$ . AR models do not make any independence assumptions but use the default factorization of  $p(x) := p(x) = \prod_{i=1}^n p(x_i | x_{\leq i})$

In AR models, these factors are  $= p(x_1) \cdot p(x_2 | x_1) \cdot p(x_3 | x_2, x_1) \cdots \cdot p(x_n | x_{\leq n})$  parameterized using a neural network and these parameters are learnt to eventually learn these factors which are used to construct the joint distribution.

The network is set up ~~such that~~ such that, at the output layer, we predict the  $n$  factors/conditional probability distributions. The  $i^{th}$  neuron in the output layer gives  $p(x_i=1 | x_{\leq i})$ . But the network should be constructed such that  $i^{th}$  output should only be connected to the previous  $i-1$  inputs. Since ~~it~~ it only depends on these  $i-1$  variables.

In NADE (Neural Autoregressive Density Estimator), a hidden representation is ~~calculated~~ calculated using only the  $i-1$  inputs, then this hidden representation is used to calculate  $p(x_i | x_{\leq i})$ . So at the  $i^{th}$  step :-

$$h_i = \sigma(w_{\leq i}^T x_{\leq i} + b) \quad \text{where } h_i \in \mathbb{R}^d \text{ and } w \in \mathbb{R}^{dn} \\ b \in \mathbb{R}^d$$

So only first  $i-1$  columns are used to calculate  $h_i$ . Then  
 $y_i = p(x_i | x_{\neq i})$  where  $[y_i = \sum_{j=1}^d h_j + c_i]$   $\forall i \in \mathbb{N}$

Notice that for calculating  $h_i$ , no input is seen, hence we make  $h_i$  as a learnable parameter of the network. So in total there are :-  $d \times n$  (for  $W$ ) +  $d$  (for  $b$ ) +  $n \times d$  (for  $V_{i,s}$ )  
 $+ n$  (for  $c_i$ )

This quantity is linear w.r.t  $n$ , instead of exponential. This is because we use weight tying. We use the same parameters  $W$  and  $b$  at each step. In practice, we mask the input appropriately to select the correct columns of  $W$ .

$\Rightarrow$  We can use the cross entropy loss for each node at the output as the loss for training. Since each output node predicts a distribution of a binary variable we can calculate its cross entropy loss and take the sum of all the losses coming from each node.

$\Rightarrow$  Notice that these models are not meant for abstraction by design. This is because we don't get a hidden representation for the whole input together from the models. We can of course devise a way to merge the individual hidden representation to calculate an abstraction.

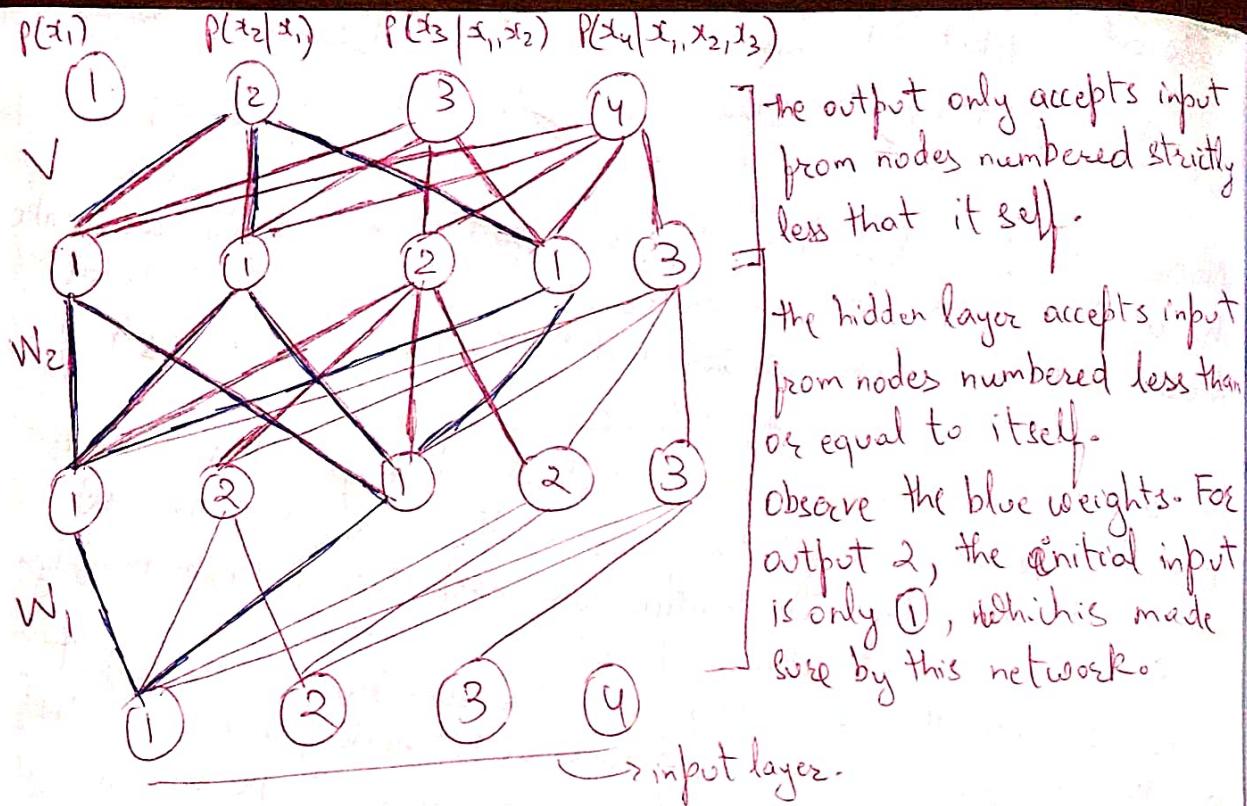
$\Rightarrow$  For generation, first the model gives  $p(x_1)$ , using this Bernoulli dist, we can sample  $x_1$  and use the sampled value to calculate  $h_2$  and then  $y_2 (p(x_2 | x_1))$ . Then  $x_2$  can be sampled and then  $x_1$  and  $x_2$  are used to calculate  $y_3$ . This process is repeated for all  $n$  vars. This process is sequential and slow. However it can be sped up by

reusing calculations made in the previous step:-

$$x_i = W_{i-1} x_{i-1} + w_{i-1} x_{i-1} \quad W_{i-1} \text{ is the } i-1^{\text{th}} \text{ col of } W \text{ and } w_{i-1} \text{ is a scalar.}$$

already calculated in previous step

Another autoregressive Model is MADE (Masked Autoencoder Density Estimator) which also considers the natural factorization. It uses the same architecture as an autoencoder, but introduces masks in the hidden layer computations which makes sure that the final  $i^{\text{th}}$  output depends on the previous  $i-1$  inputs. We start by assuming some ordering on the inputs and just number them from 1 to  $n$ .



$\Rightarrow$  Each node is randomly assigned a number between 1 to  $n-1$    
 $\text{hidden}$    
 Only for understanding, in practice,   
 only masking is done.

Let  $w_2$  be masked

$$w_2 = \begin{bmatrix} w_{11}^{(2)} & w_{12}^{(2)} & \dots & w_{15}^{(2)} \\ w_{21}^{(2)} & w_{22}^{(2)} & \dots & w_{25}^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ w_{51}^{(2)} & w_{52}^{(2)} & \dots & w_{55}^{(2)} \end{bmatrix} \quad 5 \times 5$$

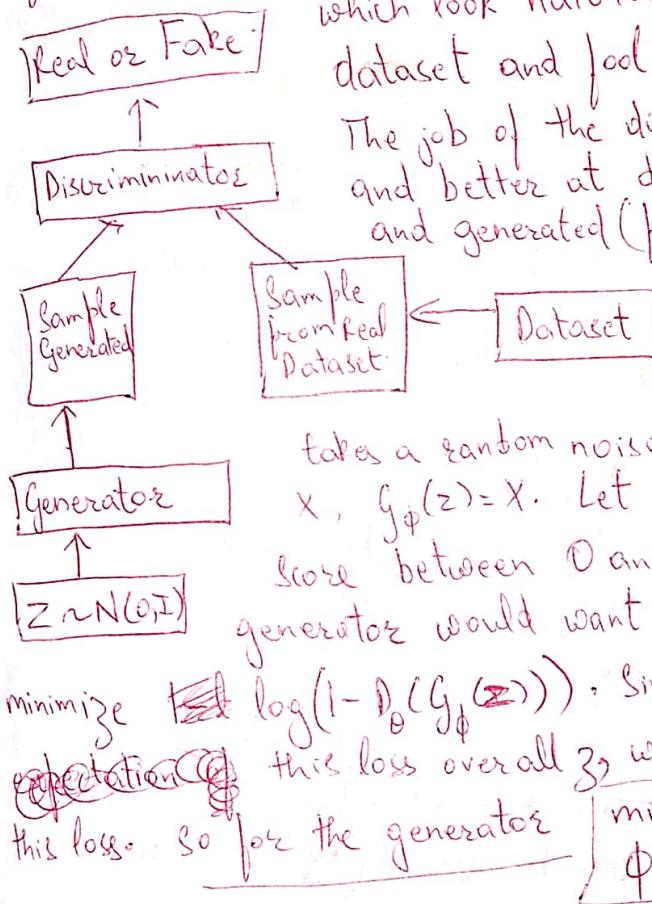
$$\text{mask} = \begin{bmatrix} 1 & 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\text{input from first hidden layer}$$

Since the first node of 2<sup>nd</sup> hidden layer depends only on 1<sup>st</sup> and 3<sup>rd</sup> node of 1<sup>st</sup> hidden layer, hence the mask contains [1 0 1 0 0] as the first row. The 3<sup>rd</sup> row of 2<sup>nd</sup> hidden layer depends only on first 4 nodes of 1<sup>st</sup> hidden layer. So the 3<sup>rd</sup> row of the mask is [1 1 1 1 0].

The loss function is same sum of cross entropy loss, and it follows the same autoregressive generation process, as in NADE

Generative Adversarial Networks :- Instead of learning a joint distribution over the data & random variables, GANs take a different approach to this problem where the idea is to sample from a simple tractable distribution ; say  $z \sim N(0, I)$  ; and then learn a complex transformation from this to the training distribution. The training of the network is like a two player game.



The job of the generator is to produce images which look natural as if sampled from the real dataset and fool the discriminator.

The job of the discriminator is to get better and better at distinguishing between true images and generated (fake) images.

Let  $G_\phi$  be the generator and  $D_\theta$  be the discriminator.  $G_\phi$

takes a random noise vector  $z \sim N(0, I)$  and produces  $x$ ,  $G_\phi(z) = x$ . Let the discriminator output a score between 0 and 1. For a given  $z$ , the generator would want to maximize  $\log D_\theta(G_\phi(z))$  or minimize  ~~$\log(1 - D_\theta(G_\phi(z)))$~~ . Since we want to minimize this expectation of this loss overall  $z$ , we minimize the expectation of this loss. So for the generator

$$\min_{\phi} \mathbb{E}_{z \sim N(0, I)} [\log D_\theta(G_\phi(z))]$$

The discriminator on the other hand has to assign a high score to real images and low score to generated images. It should do for all possible real images and all possible fake images.

$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_\theta(x)] + \mathbb{E}_{z \sim N(0, I)} [\log(1 - D_\theta(G_\phi(z)))]$$

Combining both

the objectives of the generator and discriminator :-

$$\min_{\phi} \max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log D_\theta(x)] + \mathbb{E}_{z \sim N(0, I)} [\log(1 - D_\theta(G_\phi(z)))]$$

So the overall training proceeds by

alternating between these two steps :-

i) Gradient Ascent on Discriminator

ii) Gradient Descent on generator:

In practice, the above generator objective does not work well and we use a different objective:-  
we minimize  $-\log D_\phi(G_\phi(z))$ .

GAN Training :-

for number of iterations:-

for  $K$  steps

sample  $m$  noise samples  $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$

sample  $m$  examples  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$  from data.

update the discriminator by ascending its gradient:-

$$\nabla_{\theta} \perp \sum_{i=1}^m [\log D_\theta(x^{(i)}) + \log(1 - D_\theta(G_\phi(z^{(i)})))]$$

Sample  $m$  noise samples  $\{z^{(1)}, z^{(2)}, \dots, z^{(m)}\}$

~~Sample  $m$  examples  $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$  from data.~~

update the generator by ascending its gradient:-

$$\nabla_{\phi} \perp \sum_{i=1}^m \log D_\theta(G_\phi(z^{(i)}))$$

Architecture :- For discriminator any CNN based classifier with 1 class at the output can be used (eg: VGG, ResNet).

For the generator, we have to go from a random vector to an image, called transpose convolution.

Pooling layers are replaced by strided convolutions (discriminator) and fractional-strided convolutions (generator). ReLU is used in generator except at output where tanh is used as activation function. Use Leaky ReLU in the discriminator for all layers.

It can be formally proved that when solve the above objective, the distribution learned by the generator over the variables is same as the ~~data~~ actual distribution.

Theorem: - The global minimum of the virtual ~~Perceptron~~ training criterion  $\mathcal{L}(G) = \min_{\theta} V(G, \theta)$  is achieved if and only if  $P_g = P_{\text{data}}$ .  
 So even though GANs do not model  $P(X)$ , they implicitly learn  $P(X)$  through an indirect modified objective.

Summary:

	RBM's	VAE's	AE models	GANs
Abstraction	Yes	Yes	No	No
Generation	Yes	Yes	Yes	Yes
Compute $P(X)$	Intractable	Intractable	Tractable	No
Sampling	Gibbs Sampling	Fast	Slow	Fast
Type of GM	Undirected	Directed	Directed	Directed
Loss	KL Divergence	KL Divergence	KL Divergence	Jensen-Shannon Divergence
Assumptions	$X$ independent given $Z$	$X$ independent given $Z$	None	Not applicable since $P(X)$ is not computed.
Samples	Bad	OK	Good	Good (best)

Recent works combine these models e.g., Adversarial Autoencoders, PixelGAN Autoencoders, etc.