The top layer $V_T$ is called the output layer.

In simple prediction problems, the output layer contains a single neuron whose output is the output of the network.

$h(x) \in \mathbb{R}$

We refer to $T$ as the number of layers/depth of the network. The size of the network is $|V|$. #nodes
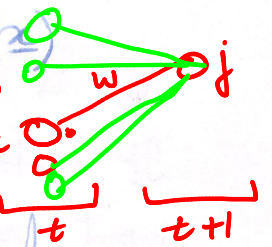
The width of the network is $\max_t |V_t|$

Suppose we have calculated the outputs of the neurons at $V_t$. Layer $t+1$ as then computes the activation value for every neuron $j$ as

E: edges of the graph

layer, neuron id.

Affine transform (layer #, node id)

scalar value

$$a_{t+1, j}(\underline{x}) = \sum_r \omega(v_{t,r}, v_{t+1,j}) O_{t,r}(\underline{x})$$

edge weights

$r: (v_{t,r}, v_{t+1,j}) \in E$

activation computed by node $j$ in layer $t+1$

o/p of the prev layer



Applying the non-linearity to the activation value:

$$O_{t+1,j}(\underline{x}) = \sigma(a_{t+1,j}(\underline{x}))$$



## Forward Propagation



Matrix $A$ is formulated using layer weights $\underline{W}_{t-1}$

Activations $\underline{a}_t = A \, O_{t+1} = B \, \underline{w}_{t-1}$ where matrix $B$ is formulated using $O_{t-1}$

# Computations in a layered Network

## Report the computation at (every node)



$\underline{O}_{t-1}$  $O_{t}$  vector

weighted sum input $b$ → activation

sigmoid/ReLU

$a_{t,j}$ → $\sigma$ Non linearity → $O_{t,j}$ scalar value.
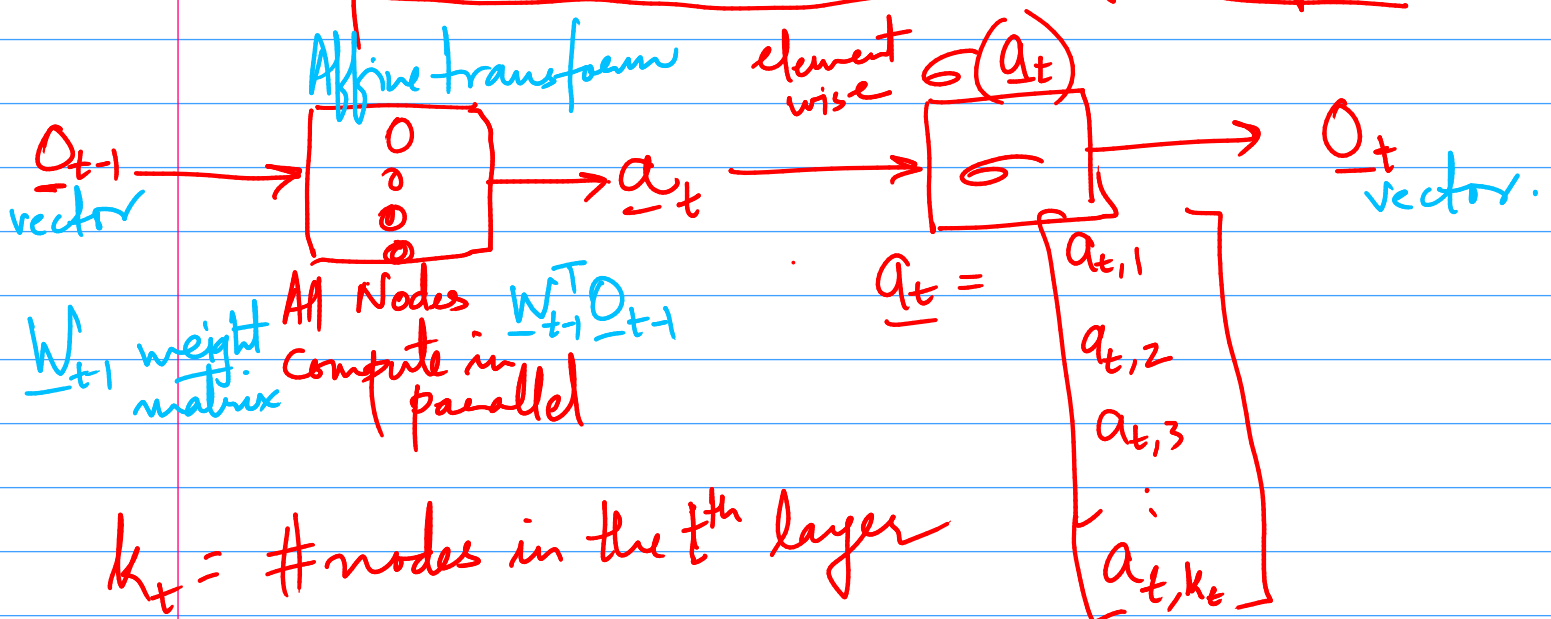
weighted sum

$j^{th}$ node of layer $t$

$k_{t-1}$ neurons

$$\underline{O}_{t-1} \equiv \begin{bmatrix} O_{t-1,1} \\ O_{t-1,2} \\ \vdots \\ O_{t-1,k_{t-1}} \end{bmatrix} \qquad \underline{w} = \begin{bmatrix} w(v_{t-1,1}, v_{t,j}) \\ w(v_{t-1,2}, v_{t,j}) \\ \vdots \end{bmatrix}$$

$$a_{t,j} = \underbrace{\underline{w}^{T} \underline{O}_{t-1}}_{\text{inner product}} = \underline{O}_{t-1}^{T} \underline{w}$$

$\mathbb{R}$. scalar output

$$O_{t,j} \equiv \sigma(a_{t,j})$$

## Computations carried out by a layer



Affine transform

element wise $\sigma(\underline{a}_t)$

$\underline{O}_{t-1}$ vector →  → $\underline{a}_t$ →  $\sigma$ → $\underline{O}_t$ vector.

$\underline{W}_{t-1}$ weight matrix.

All Nodes compute in parallel $\underline{W}_{t-1}^{T} \underline{O}_{t-1}$

$$\underline{a}_t = \begin{bmatrix} a_{t,1} \\ a_{t,2} \\ a_{t,3} \\ \vdots \\ a_{t,k_t} \end{bmatrix}$$

$k_t = \#$ nodes in the $t^{th}$ layer

# Architecture $(V, E)$

layers.



**Dense Connection** (green) · **Sparse Connection** (red) · FF · skip connection · $(V, E)$

$w$: learnable parameters

$V_{0,1}$

$n.$   $V_{on}$

$V_0$

Bais-Variance trade off.

$\gtrless ? \longrightarrow$  Architecture (edges)

$\begin{cases} \text{\# layers} \\ \text{\# nodes in each layer} \\ \text{Connections between nodes} \end{cases}$

A → σ → B → σ → C → σ → output

$$W = ABC$$

A B C →
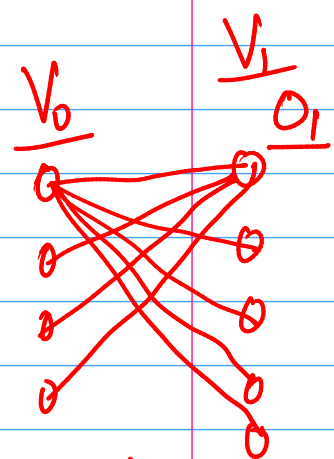
provided $\sigma(\underline{a}) = I$ identity

Behaves as a linear transformation

Universal Approximators ⎤ Non linearity
                        ⎦ is important

$\mathcal{H}_{V,E\beta}$

$h \in \mathcal{H}$ which can correctly model it

$V_0$   $V_1$
        $O_1$   $O_2$   $O_{t+1}$   $O_t$

O       O       O       O       O
O       O       O               O
O       O       O       O       O
O       O       O       O       O
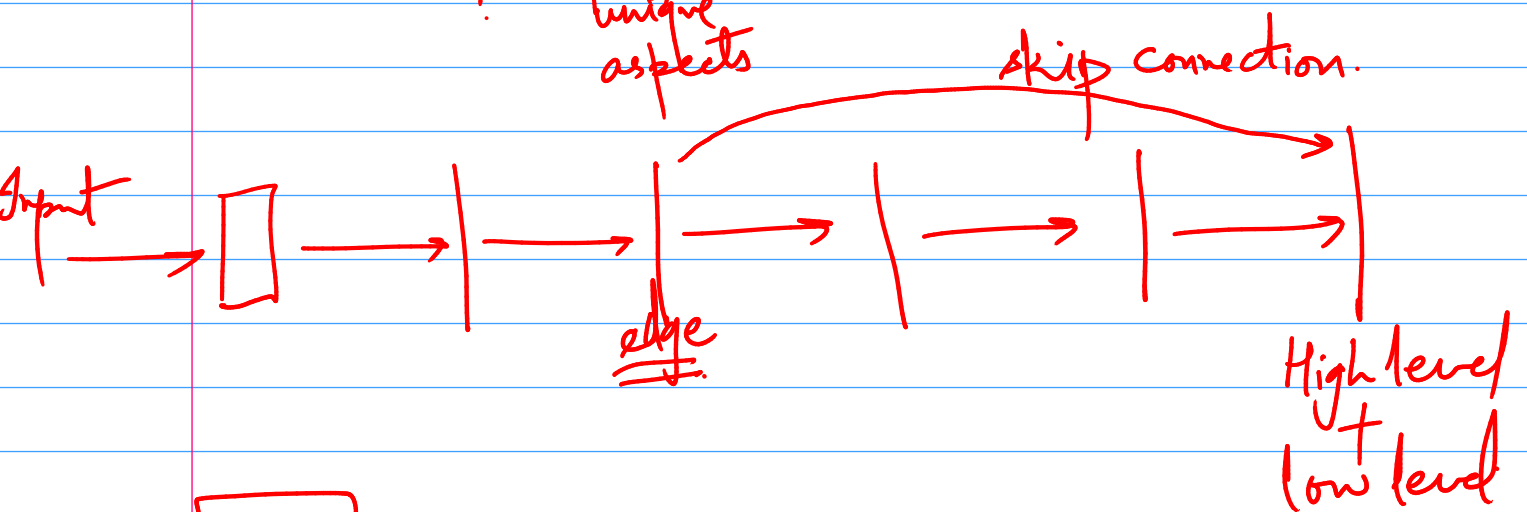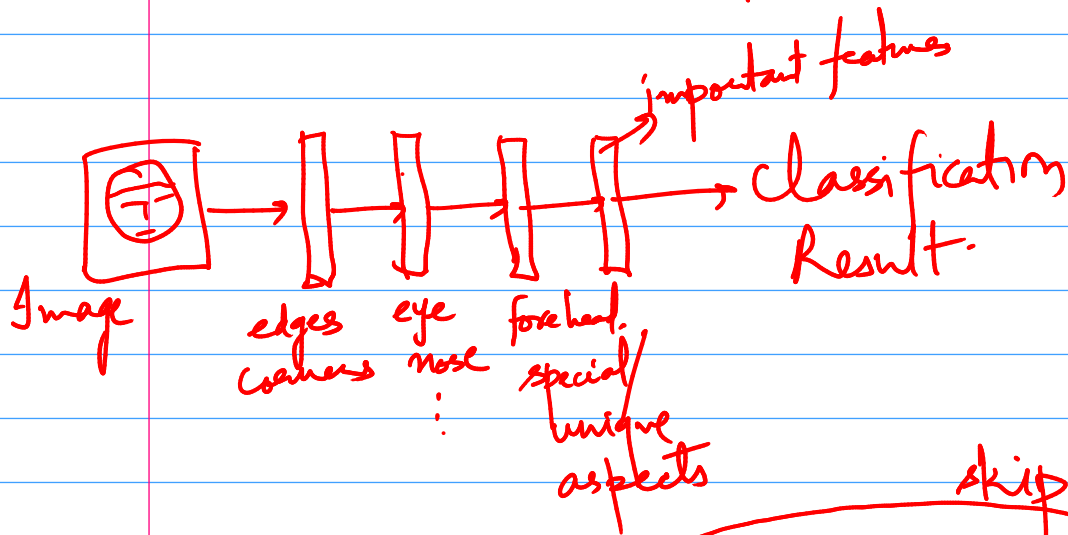        O       O       O       O
        O       O       ⟶ feed forward network.

⎵
input layer

no computations

Features of higher level of abstraction are computed by the deeper layers

low level features are composed to form high level features

important features



classification Result.

Image

edges eye forehead.
corners nose special
: unique
aspects
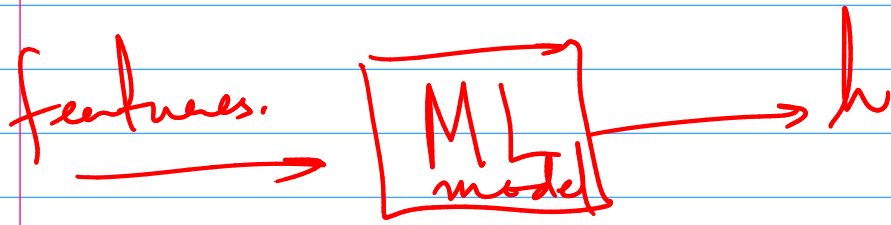
skip connection.

Input

edge

High level
+
low level

SGD

Backpropagation computes the gradients of the loss w.r.t. parameter of each & every layer.

Weight update using gradients

Learning features automatically.

features. $\longrightarrow$ [ ML model ] $\longrightarrow h$

Engineered / Handcrafted designed. | Learned.

Representation

Learning

$\underline{x}$ $\begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix} \longrightarrow \begin{bmatrix} - \\ - \\ - \end{bmatrix}$