# Forward Propagation



$$a_t = \underline{A}\, \underline{O}_{t-1} = \underline{B}\, \underline{w}_{t-1}$$

$W_{t-1}$   $k_{t-1} \times k_t$

$\underline{w}_{t-1} = $   $k_{t-1} k_t \times 1$
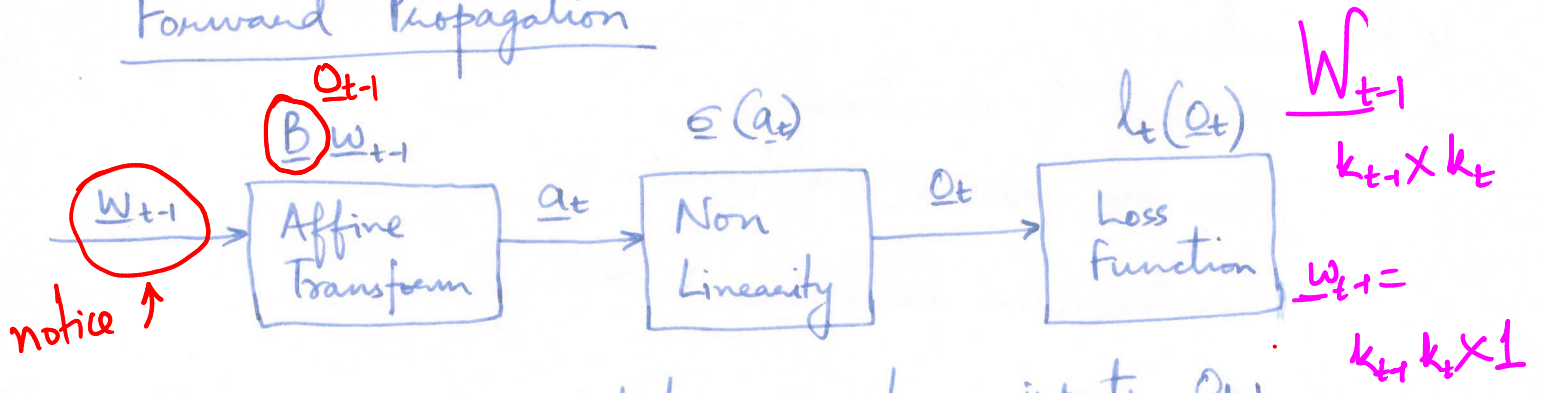
Matrix $\underline{B}$ is formulated using layer inputs $\underline{O}_{t-1}$

Activations $\quad a_t = \underline{A}\, \underline{O}_{t-1} = \underline{B}\, \underline{w}_{t-1}$   → vector $\underline{w}_{t-1}$

$$\underline{W}_{t-1}\underline{O}_{t-1} = \underline{B}\, \underline{w}_{t-1}$$
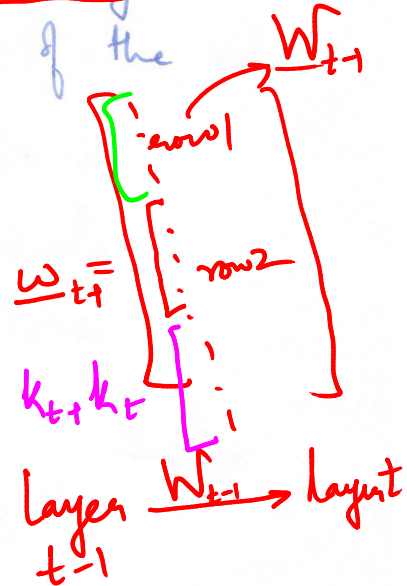
vector $\underline{w}_{t-1}$ is the column vector obtained by concatenating the rows of $\underline{W}_{t-1}$ (matrix) and then taking the transpose of the resulting long row vector. $\quad k_{t-1} k_t$.

$$\underline{B} \equiv \underline{O}_{t-1} = k_t \begin{bmatrix} \underline{O}_{t-1}^T & 0^z & \cdots \to \text{zeros} & \cdots & 0^z \\ 0^z & \underline{O}_{t-1}^T & & & 0^z \\ \vdots & & & & \\ 0^z & 0^z & & & \underline{O}_{t-1}^T \end{bmatrix}$$

$$B \equiv \begin{bmatrix} \cdots & 0 & 0 & 0 \\ 0 & \cdots & 0 & 0 \\ 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & 0 & 0 & \cdots \end{bmatrix}$$

$$\underline{w}_{t-1} = \begin{bmatrix} \text{row 1} \\ \vdots \\ \text{row 2} \\ \vdots \\ \vdots \end{bmatrix} \quad W_{t-1}$$

$k_{t-1} k_t$

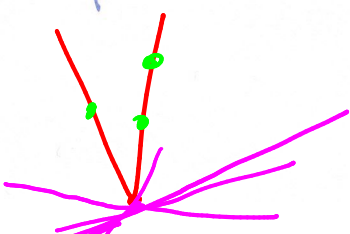Layer $\underline{W}_{t-1} \to$ layer $t$
$t-1$

If $k_t$ denotes the number of neurons in the layer $V_t$
then size of $\underline{O}_{t-1}$ is $\quad k_t \times (k_{t-1} k_t)$

vector $\underline{w}_{t-1}$ is $\quad (k_{t-1} k_t) \times 1$.

Training a Neural Network : Stochastic Gradient Descent (SGD).

Inputs to the SGD:
- Training examples $(\underline{x}, y)$
- Layered Graph $(V, E)$
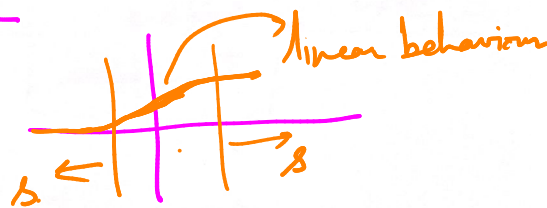- A differentiable non-linearity $\sigma$

**(H.) Hyper Parameters:**

Number of iterations $\tau$

Step size sequence $\eta_1 \, \eta_2 \cdots \eta_\tau$

Regularization parameter $\lambda > 0$

$\underline{O} = \underline{A}\,\underline{O}_{t-1}$

linear behaviour

$s \leftarrow \| \cdot \| \rightarrow s$

**Initialize:**

Choose $\underline{w}^{(1)} \in \mathbb{R}^{|E|}$ at random from a distribution such
that $\underline{w}^{(1)}$ is close to $\underline{O}$ zerovector.

$\rightarrow$ #edges

small weights ensure small activations.

for $i = 1, 2, \ldots \; (\tau)$

    sample $\underline{(x, y)} \sim D$

    calculate gradient $\underline{v_i} = $ backpropagation $(\underline{x}, \underline{y}, \underline{w}, (V,E), \sigma)$  ? ? ? ?

    update $\underline{w}^{(i+1)} = \underline{w}^{(i)} - \eta_i \left( \underline{v_i} + \lambda \underline{w}^{(i)} \right)$    $\frac{1}{2}\|\underline{w}\|^2$

               SGD update step.

Output:
    $\underline{w}^*$ is the best performing $\underline{w}^{(i)}$ on the validation set
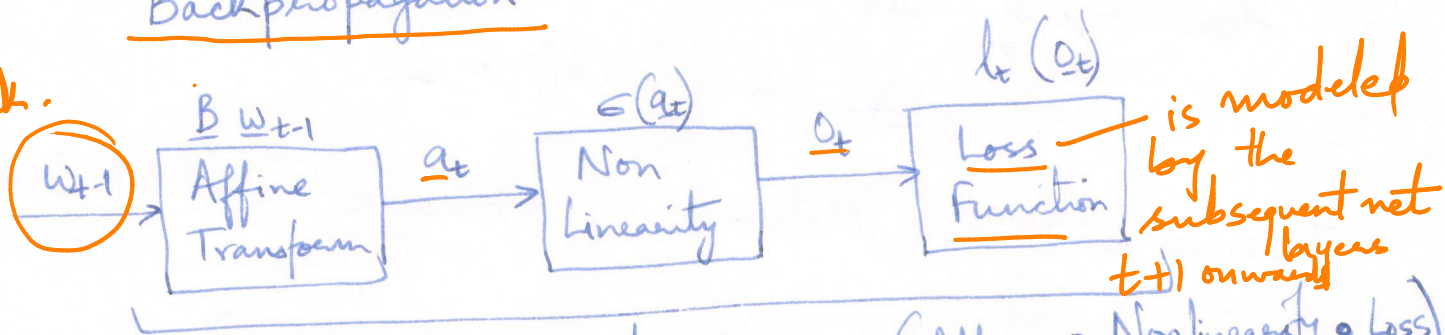
---

Computing the Gradients

Backpropagation   no space
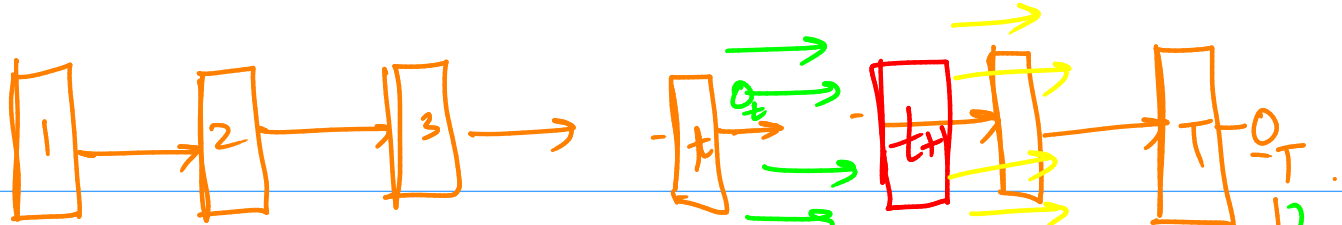
back ~~propagation~~

think.

$(\underline{w}_{t-1})$ $\xrightarrow{\underline{B}\,\underline{w}_{t-1}}$ [ Affine Transform ] $\xrightarrow{\underline{a}_t}$ [ Non Linearity ] $\xrightarrow{\underline{O}_t}$ [ Loss Function ]

                        $\sigma(\underline{a}_t)$              $l_t(\underline{O}_t)$

is modeled by the subsequent net layers $t+1$ onwards

Composite function $g_t = ($ Affine $\circ$ Nonlinearity $\circ$ Loss $)$

$g_t(\underline{w}_{t-1})$

nested function

$h\left(g\left(f(\underline{x})\right)\right)$    $f \circ g \circ h$   composition of function

$$l_t(\underline{O}_t) = l_{t+1}(\underline{O}_{t+1}) = l_{t+2}(\underline{O}_{t+2}) = l_T(\underline{O}_T) \quad l_t(\underline{O}_t) \quad \boxed{loss}$$

$$= l_T(\underline{O}_T) \; |E| \quad l_{t+1}(\underline{O}_{t+1})$$

Parameter space.

loss

loss value
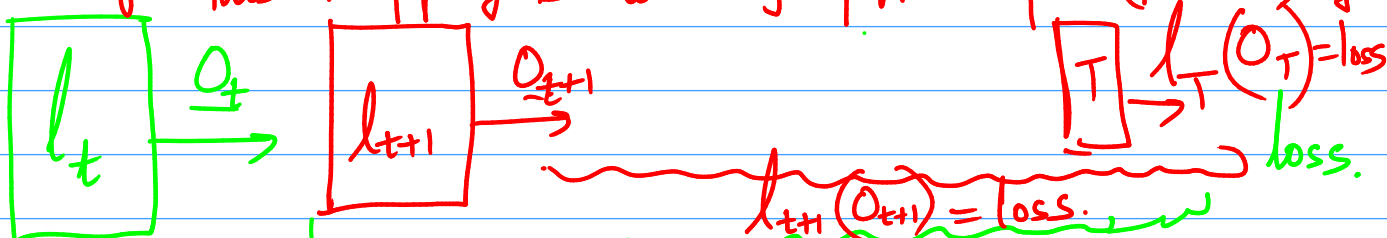
$|E|$ dim parameter space

large $\underline{w}$ $\Rightarrow$ large values for activations.

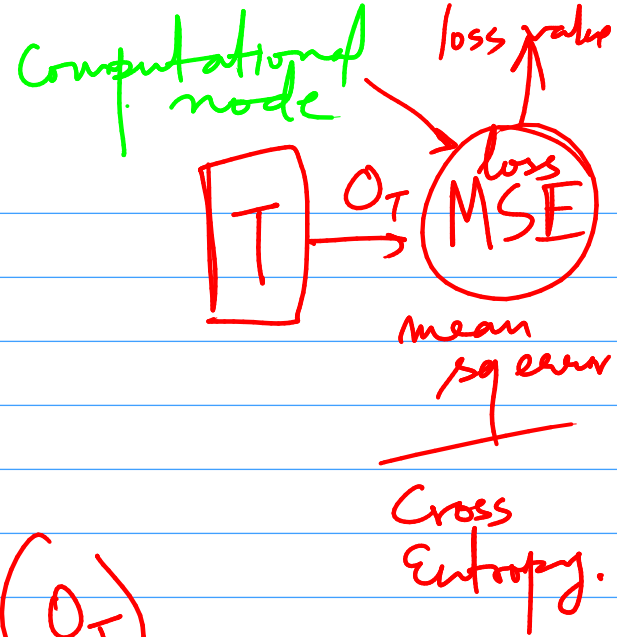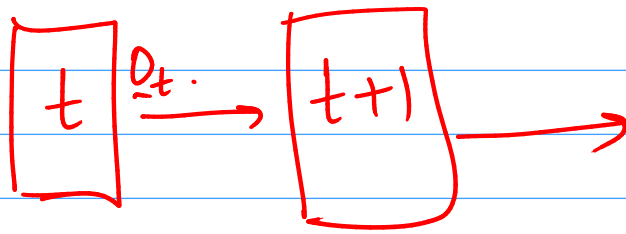If non lin are saturating $\longrightarrow$ out values are also saturated

Gradients are too small (vanishing gradient problem).

small gradients $\Rightarrow$ small updates to $\underline{w}$

$\Rightarrow$ slow learning.

outputs of each layer can be mapped to the loss val
This mapping is done by the subsequent (further) layers

$$\boxed{l_t} \xrightarrow{\underline{O}_t} \boxed{l_{t+1}} \xrightarrow{\underline{O}_{t+1}} \quad \boxed{T} \xrightarrow{} l_T(\underline{O}_T) = loss$$

$$\longrightarrow loss.$$

$$l_{t+1}(\underline{O}_{t+1}) = loss.$$

$$l_t(\underline{O}_t) = loss.$$

$$loss = l_T(\underline{O}_T) = \cdots = l_{t+1}(\underline{O}_{t+1}) = l_t(\underline{O}_t)$$

$$\boxed{t} \xrightarrow{O_t} \boxed{t+1} \longrightarrow$$

loss value

$$\boxed{T} \xrightarrow{O_T} \overset{loss}{MSE}$$

mean sq error

———

Cross Entropy.

$$loss = MS\dot{E}(O_T)$$

$$MSE(O_t)$$

$$MSE\left(\cdot \underset{t+1 \text{ onwards}}{layers} \text{ that process } O_t \text{ & map } \right)$$
$$it \quad to \quad O_T$$

$$MSE\left(Subsequent\ layers\ after\ layer\ (O_t)\right)$$

$$loss = MSE\left(\underline{l_t}\ (O_t)\right)$$

$$= MSE\left(l_{t+1}\ (O_{t+1})\right)$$

$$= \underbrace{\frac{MSE(O_T)}{(y_T - O_T)}}_{2}$$