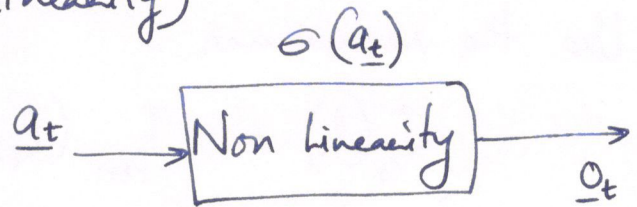


Sigmoid Function (Non-linearity)

$$\underline{o}_{ti} = \frac{1}{1 + \exp(-\underline{a}_{ti})}$$



$$\sigma'(\underline{a}_{ti}) = \frac{1}{(1 + \exp(-\underline{a}_{ti})) (1 + \exp(\underline{a}_{ti}))}$$

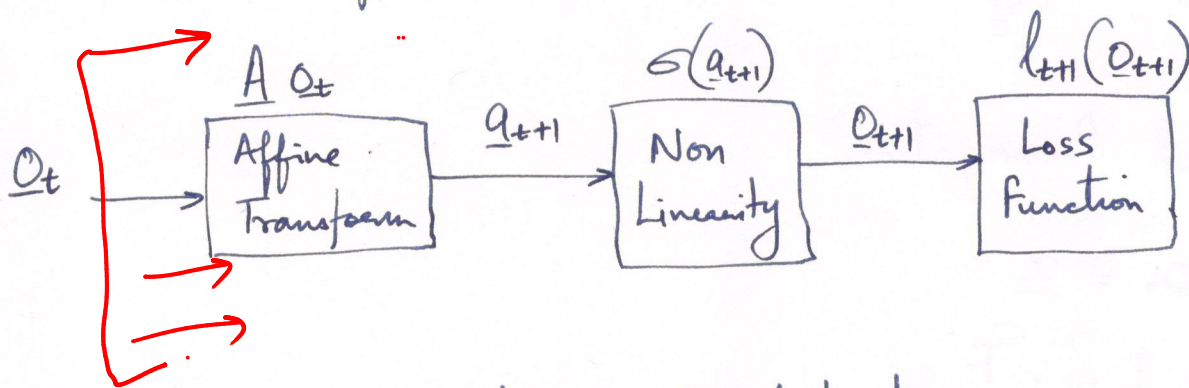
i^{th} diagonal entry of $J_{\underline{a}_t}(\underline{o}_t)$

$J_{\underline{a}_t}(\underline{o}_t)$ is a diagonal matrix denoted as $\text{diag}(\sigma'(\underline{a}_t))$
elementwise nonlinearity

$$l_t(\underline{o}_t) = l_{t+1}(\underline{o}_t) = l_{t+1}(\sigma(\underline{w}_t \underline{o}_t))$$

$$J_{\underline{a}_t}(l_t) = \underbrace{J_{\underline{a}_{t+1}}(l_{t+1})}_{\text{Loss from further network}} \underbrace{J_{\underline{a}_{t+1}}(\underline{o}_{t+1})}_{\text{Non linearity } \sigma} \underbrace{J_{\underline{o}_t}(\underline{a}_{t+1})}_{\text{Affine}}$$

recursive



Computing $J_{\underline{o}_T}(l_T)$ for the last layer

Final loss $l_T = \frac{1}{2} \|\underline{o}_T - \underline{y}\|^2$ \underline{y} : ground truth target

$J_{\underline{o}_T}(l_T) = \underline{o}_T - \underline{y}$ vector.

Use the recurrence

the recurrence

$$\gamma_t \circ J_q(l_t) = \underbrace{J_{\underline{a}_{t+1}}(l_{t+1})}_{\gamma'} \underbrace{J_{\underline{a}_{t+1}}(\underline{0}_{t+1})}_{\text{affine } W.}$$

to back propagate values for $J_{\underline{O}_T}(l_T)$ starting from the last layer T and working towards the previous layers.

Backpropagation (Training examples $(\underline{x}, \underline{y})$,
layered graph (V, E))

Initialize :

denote layers of the graph $V_0 \dots V_T$ where # nodes

$$V_t = \{v_{t,1}, \dots, v_{t,n_t}\} \quad t \rightarrow n_t$$

Define $w_{t,i,j}$ as the weight of the edge $(v_{t,j}, v_{t+1,i})$

Forward Propagation:

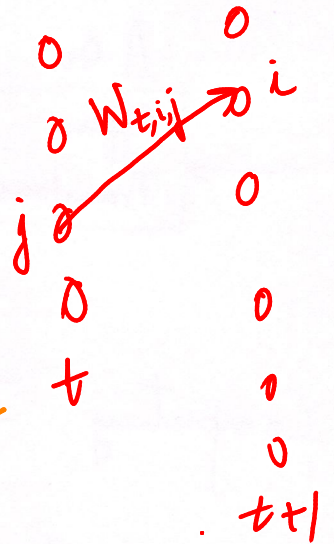
set $\underline{O}_0 = \underline{x}$ V_0 $\underline{O}_0 = \underline{x}$
 initial final layer layers
 for $t = 1, \dots, T$

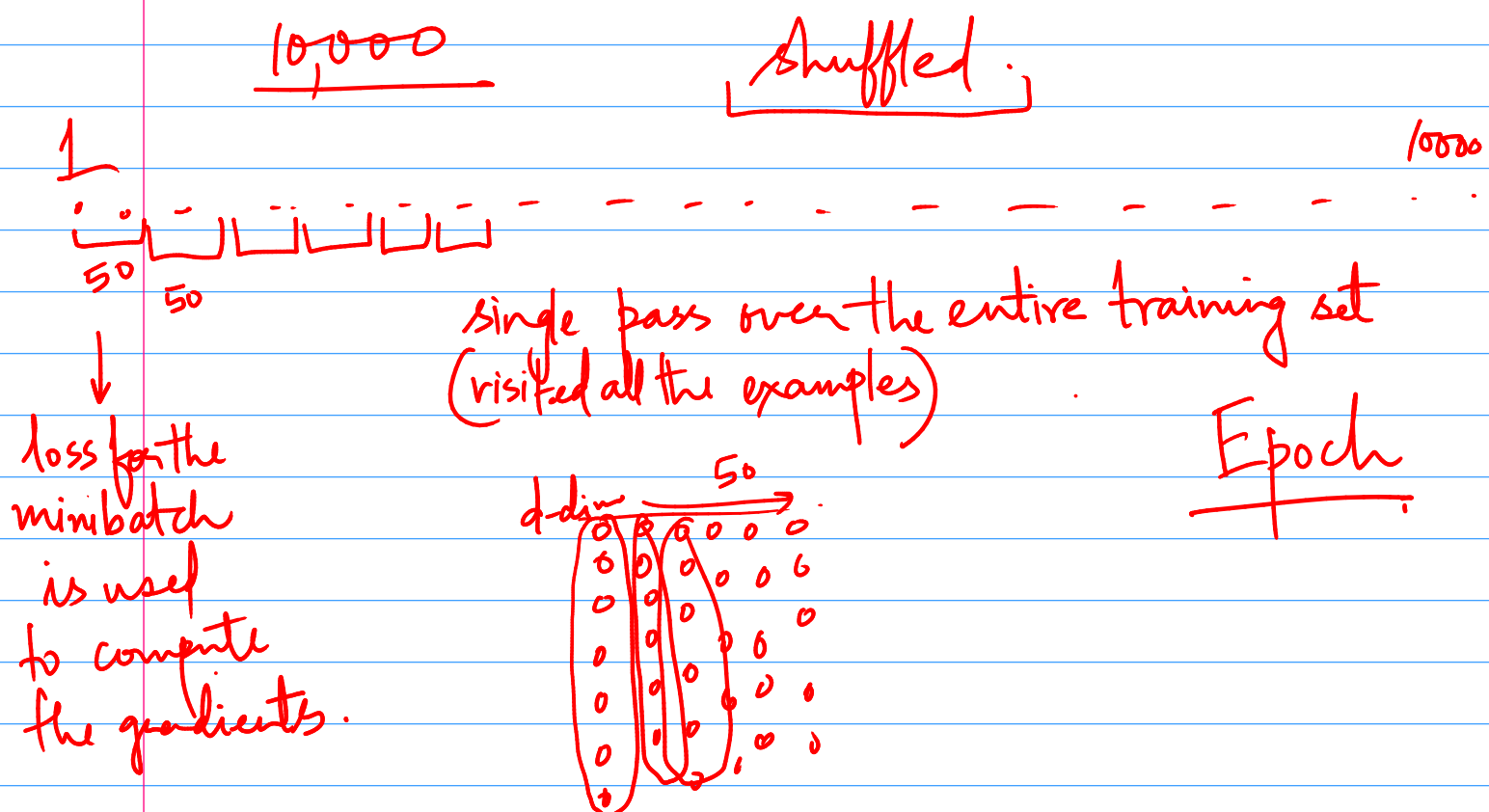
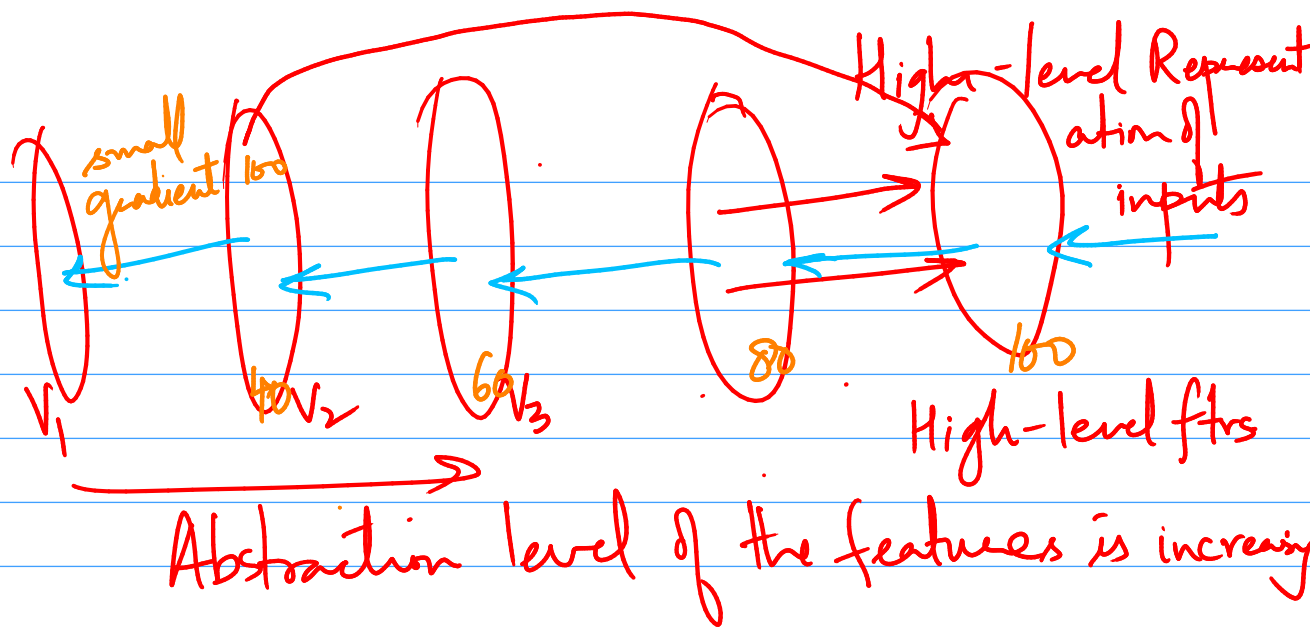
for $t = 1, \dots, T$

for $i = 1, \dots, k_t$ nodes in a layer

$$a_{t,i} = \sum_{j=1}^{k_{t-1}} \underbrace{O_{t-1,j}}_{\text{prev}} \underbrace{W_{t-1,i,j}}_{\text{weights}} \quad \text{Affine}$$

$$O_{t,i} = \underbrace{\sigma(a_{t,i})}_{\text{non-linearity}}$$





Backward Propagation :

gradient of the final loss

set $\delta_T = (O_T - y)$

where $\delta_T \equiv \underline{J_{O_T}(l_T)}$

δ_T

δ_{T-1}

δ_{T-2}

for $t = \underline{T-1}, \underline{T-2}, \dots, \underline{1}$

① for $i = 1, \dots, n_t$

$\delta_{t,i} = \sum_{j=1}^{n_{t+1}} W_{t,j,i} \delta_{t+1,j} \sigma'(a_{t+1,j})$

incoming gradients from outgoing links.

② for each edge $(v_{t+1,j}, v_{t,i}) \in E$

set the partial derivative to $\delta_{t,i} \sigma'(a_{t,i})$

$\delta_{t,i}$

used to update the weight

SGD

Loss Functions

used for Neural Network Training

Regression Tasks

Use the squared error loss

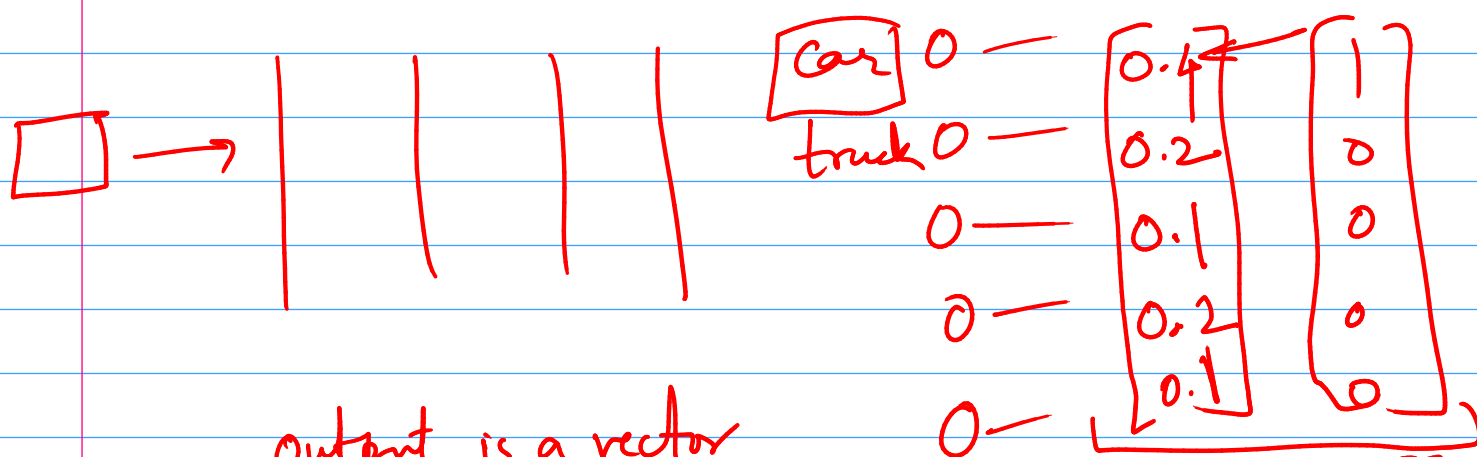
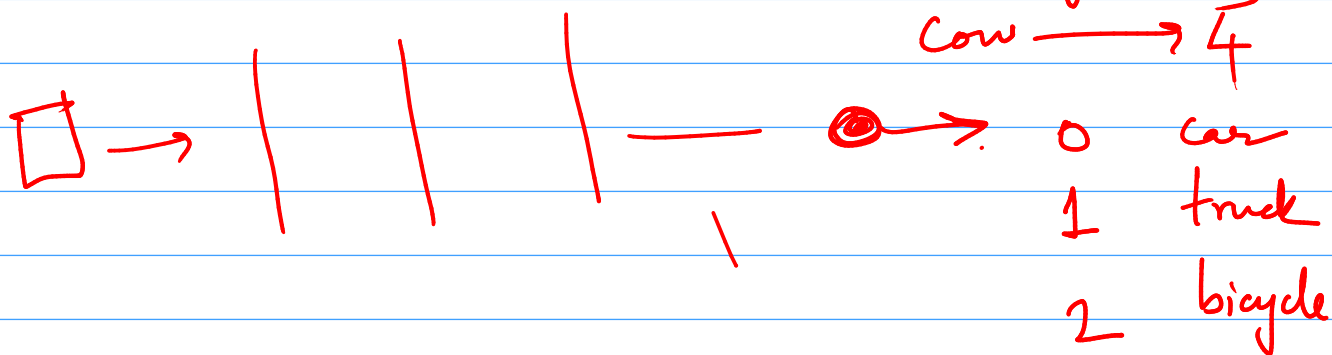
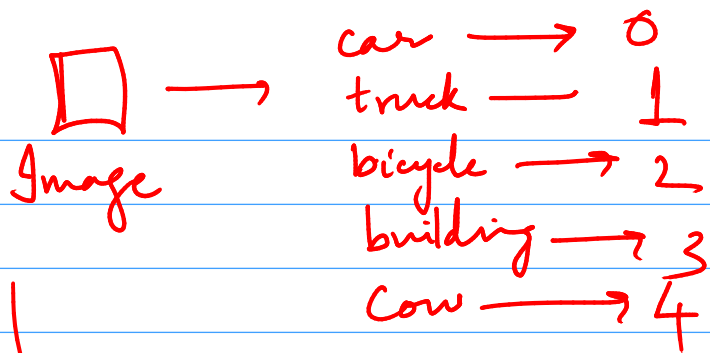
$$MSE_Loss(w) = -\frac{1}{2} \mathbb{E}_{x,y \sim S} \| \underset{\text{target}}{y} - \underset{\text{o/p of neural net}}{f(x; \theta)} \|^2$$

Classification tasks

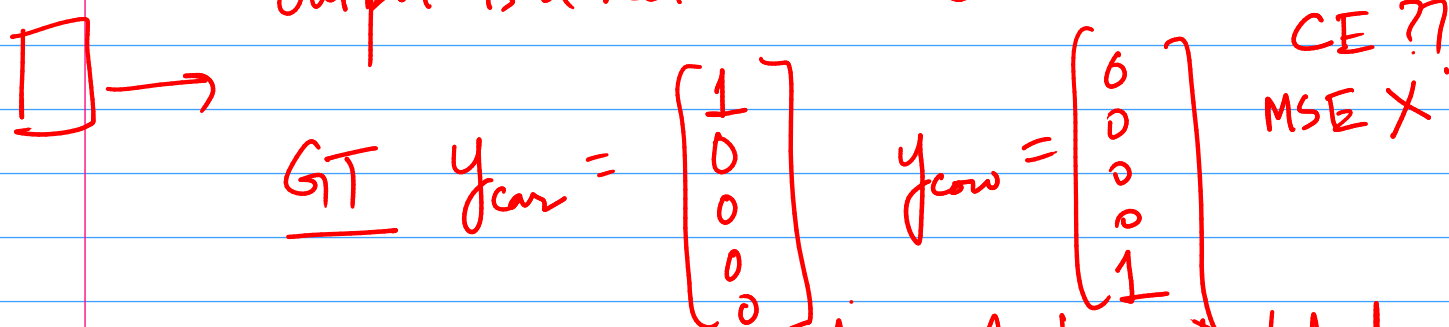
Cross entropy between the training data distribution and model distribution

$$XentropyLoss = - \mathbb{E}_{x,y \sim S} \log \left(p_{\text{model}}(y|x) \right)$$

Classification



output is a vector



1-hot encoding of target labels.

GT prob distribution of the target $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ 1-hot vector

Network produces a prob distr over the target.
Cross entropy gives distance between two prob distributions

$$CE(P, \hat{P}) = \sum_i p_i \log \frac{1}{\hat{p}_i}$$