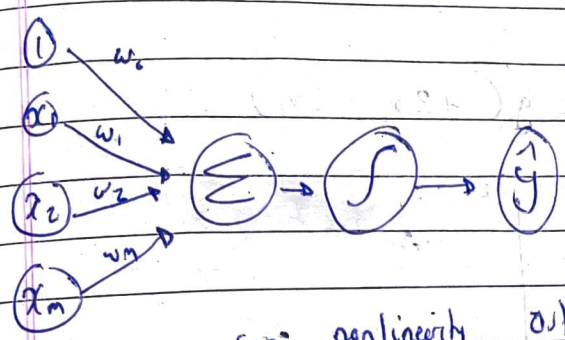


MIT 6.S191 DL

The perceptron: Forward Prop



inputs weights sum non-linearly output

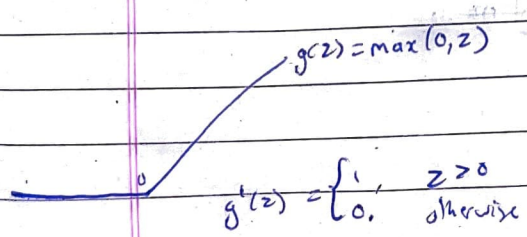
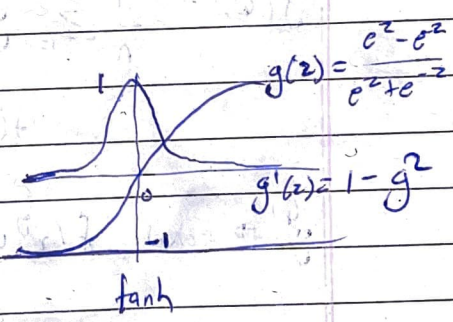
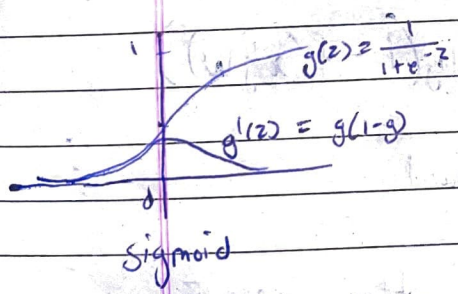
$$\hat{y} = g\left(w_0 + \sum_{i=1}^n x_i w_i\right)$$

$$= g(w_0 + X^T W)$$

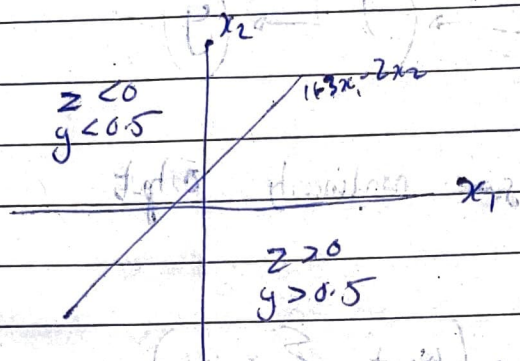
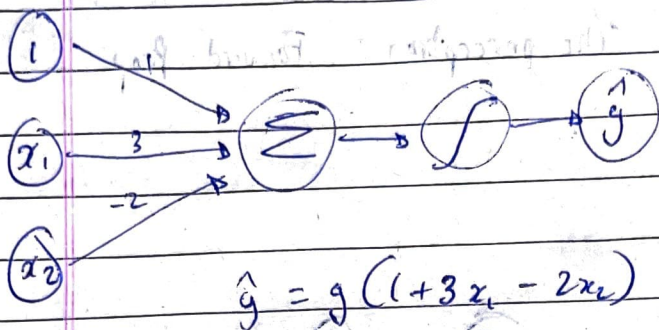
$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$W = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

Activation Functions



→ these activation funcⁿ are used to introduce non-linearity to the network



Loss functions

→ Cross entropy loss for models that have probability output. b/w 0 & 1

$$J(w) = -\frac{1}{n} \sum_{i=1}^n y^{(i)} \log(f(x^{(i)}; w)) + (1 - y^{(i)}) \log(1 - f(x^{(i)}; w))$$

$y^{(i)} \rightarrow$ actual $f(x^{(i)}; w) \rightarrow$ predicted

~~In TF, In TF, search softmax cross EL~~
~~loss = TF.reduce_mean(tf.nn.so~~

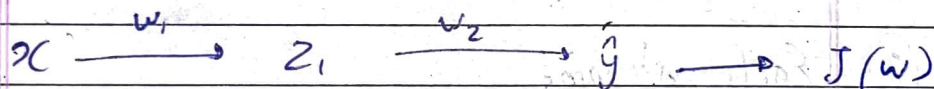
loss optimisation

$w^* \rightarrow w$ that minimises $J(w)$

Gradient descent

1. Initialise weights randomly
2. loop till ~~converge~~ converge
3. compute gradient
4. Update weight $w \leftarrow w - \eta \frac{\partial J(w)}{\partial w}$
5. Return weights

Back prop



→ how small change in one wt affects the loss

$$\frac{\partial J(w)}{\partial w_2} = \frac{\partial J(w)}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial w_2}$$

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

→ Stable learning rates converge smoothly and avoid local min

Adaptive Learning rates

Gradient descent algo

- SGD
- Adam
- Adadelta
- Adagrad
- RMS Prop

Regularisation

During training

① randomly set some activⁿ to 0

↓
Dropout

- Forces network to not rely on any 1 node
- Drops different every iteration

② Early stopping

