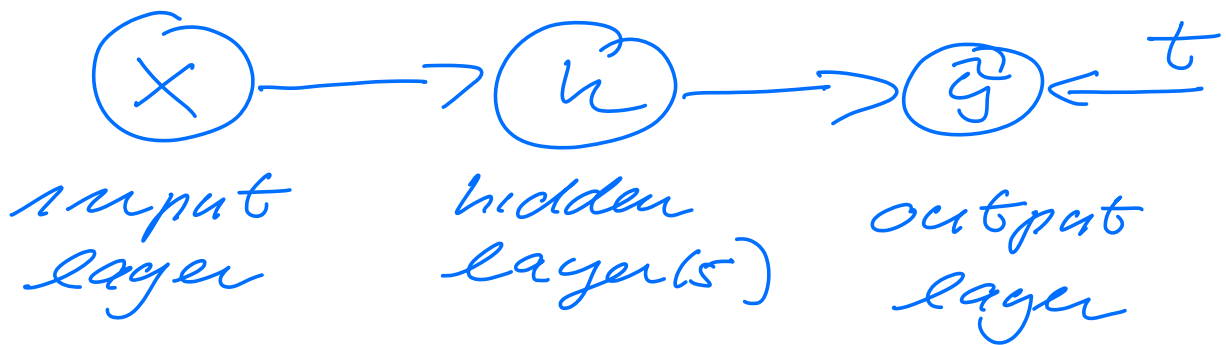


FYS-STK 4155, OCT 7, 2022

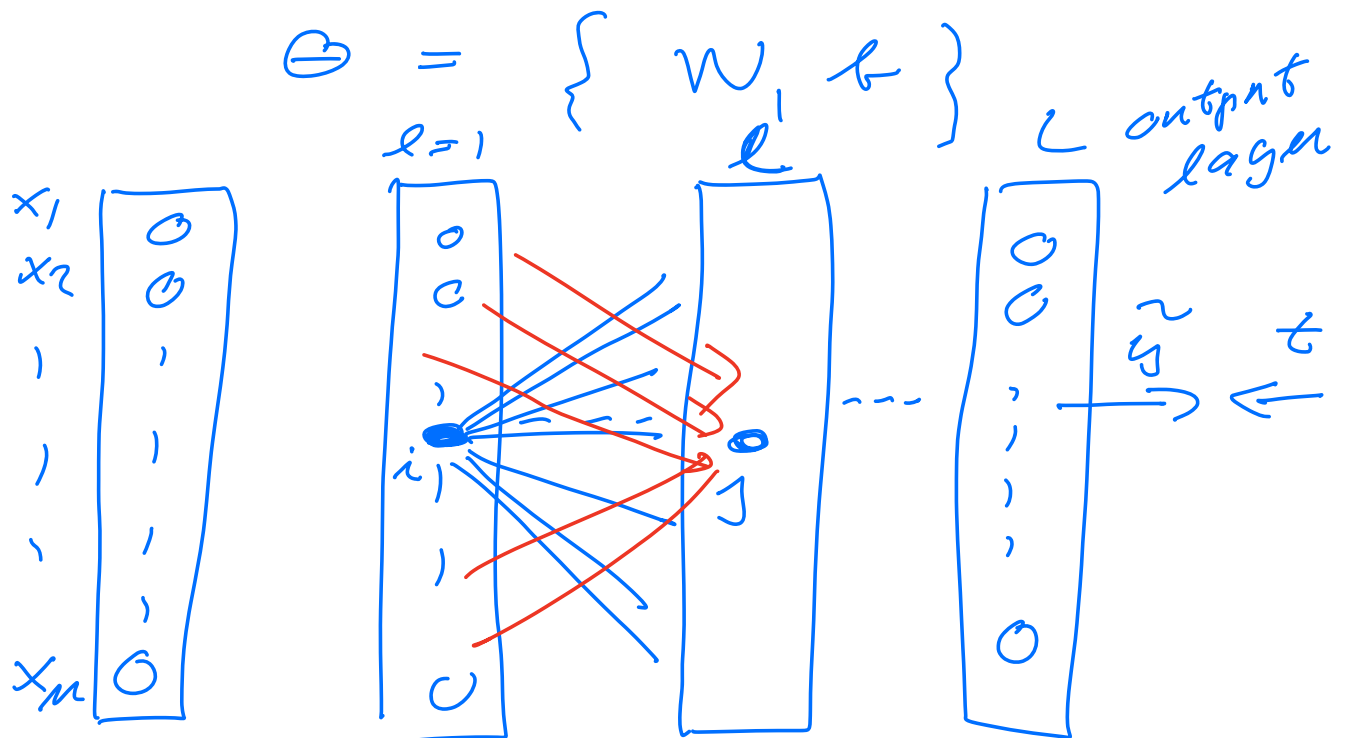
Neural Networks



Multi perceptron (MLP)

Feed Forward NN (FFNN)

our parameters



N_l N_l

input to node j in layer
(hidden) - l - from layer
- $l-1$ -

$$z_j^l = \sum_{i=1}^{N_{l-1}} w_{ij}^l a_i^{l-1} + b_j^l$$

a_i^{l-1} = output from
node i in layer
 $l-1$

Cost/Loss function (regression)

$$C(\theta) = \frac{1}{2} \sum_{i=1}^n (a_i^L - t_i)^2$$

$$a_j^l = \sigma^l(z_j^l)$$

activation function

$$z^l = (W^l)^T a^{l-1} + b^l$$

Popular activation function

Sigmoid

$$\begin{aligned}\sigma^l(z_j^l) &= \frac{1}{1 + e^{-z_j^l}} \\ &= a_j^l\end{aligned}$$

Basic ingredients in FFNN

- Feed forward stage
- Back propagation stage to train Θ
 - Expression for gradients of C as function of the various layers.
- Gradient optimizer (SGD, momentum)

Adaptive ...)

- Repeat with Feed forward and back propagation to update Θ till our cost function reaches a converged value.

Back propagation algo

$$z_j^l = \sum_{i=1}^{N_{l-1}} w_{ij}^l a_i^{l-1} + b_j^l$$

$$a_i^{l-1} = \sigma(z_i^{l-1})$$

$$\frac{\partial z_j^l}{\partial w_{ij}^l} = a_i^{l-1}$$

$$\frac{\partial z_j^l}{\partial a_i^{l-1}} = w_{ij}^l$$

$$\frac{\partial a_j^l}{\partial z_j^l} = ?$$

$$a_j^l = \sigma^l(z_j^l) = \frac{1}{1 + e^{-z_j^l}}$$

$$\frac{\partial a_j^l}{\partial z_j^l} = \underbrace{\sigma^l(z_j^l) (1 - \sigma^l(z_j^l))}_{\text{changes with activation function}}$$

$$= a_j^l (1 - a_j^l)$$

$l = L$ (output layer)

$$C(\theta) = \frac{1}{2} \sum_i (a_i^L - t_i)^2$$

$$\frac{\partial C(\theta)}{\partial w_{jk}^L} = (a_i^L - t_i) \frac{\partial a_j^L}{\partial w_{jk}^L}$$

$$\begin{aligned}
 \frac{\partial a_j^L}{\partial w_{jk}^L} &= \frac{\partial a_j^L}{\partial z_j^L} \frac{\partial z_j^L}{\partial w_{jk}^L} \\
 &= \underbrace{a_j^L (1 - a_j^L)}_{\frac{\partial \sigma(z_j^L)}{\partial z_j^L}} a_k^{L-1}
 \end{aligned}$$

$$\frac{\partial C}{\partial w_{jk}^L} = (a_j^L - t_j) \cancel{a_j^L (1 - a_j^L)} \times a_k^{L-1}$$

$$= (a_j^L - t_j) \sigma'(z_j^L) a_k^{L-1}$$

$$\delta_j^L = \sigma'(z_j^L) (a_j^L - t_j)$$

$$= \sigma'(z_j^L) \frac{\partial C}{\partial a_j^L}$$

in vector notation

$$\delta^L = \nabla^L(z^L) \odot \frac{\partial C}{\partial a^L}$$

$$\begin{bmatrix} x_1' \\ x_2' \end{bmatrix} \odot \begin{bmatrix} y_1' \\ y_2' \end{bmatrix} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \end{bmatrix}$$

$$\boxed{\frac{\partial C}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}}$$

$$\delta_j^L = \frac{\partial C}{\partial z_j^L} = \frac{\partial C}{\partial a_j^L} \frac{\partial a_j^L}{\partial z_j^L}$$

$$\delta_j^L = \frac{\partial C}{\partial b_j^L} \frac{\partial b_j^L}{\partial z_j^L} = \frac{\partial C}{\partial b_j^L}$$

$$\boxed{l = L}$$

$$\frac{\partial C}{\partial w_{jk}^L} = \delta_j^L a_k^{L-1}$$

$$\delta_j^L = \nabla^L(z_j^L) \frac{\partial C}{\partial a_j^L}$$

$$\frac{\partial C}{\partial t_j^L} = \delta_j^L$$

$$l \neq L$$

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$= \sum_k \delta_k^{l+1} \frac{\partial z_k^{l+1}}{\partial z_j^l}$$

$$z_j^{l+1} = \sum_{n=1}^{N_l} w_{nj}^{l+1} a_n^l + b_j^{l+1}$$

$$\uparrow$$

$$\Delta(z_j^l)$$

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \Delta'(z_j^l)$$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} \Delta'(z_j^l)$$

Algo

- Define Model
(architecture of NN)
 - # hidden layers
 - # nodes in layers
- initialize $\{W, b\} = \Theta$
- activation function
- hyperparameters
 - λ , l_2 or l_1 regularization
- learning rate +
learning rate scheduler
 - Adagrad
RMS Prop ...

- Gradient descent method
- Need X and t (target)

while convergence not met

1 First Feed Forward

$$\text{output} = a^L = \sigma^L(\sigma^{L-1}(\dots \sigma'(z') \dots))$$

$$z' = W^L X + b$$

2 First Back prop stage

$$l = L-1, L-2, \dots, 1$$

$$\delta_j^l = \sum_k \delta_k^{l+1} w_{kj}^{l+1} \sigma'(z_j^l)$$

$$w_{jk}^l \leftarrow w_{jk}^l - \underbrace{\eta}_{\text{learning rate}} \delta_j^l a_k^{l-1}$$

$$b_j^l \leftarrow b_j^l - \eta \delta_j^l$$

3 back to 1 and 2

repeat 3 till convergence

Problems: Vanishing gradients,

Example $L=2$, simple case, x , w and b are scalars.

$$\text{output } f(x; \theta) = \hat{y} \\ = \sigma_2(w_2 \sigma_1(w_1 x + b_1) + b_2)$$

$$\partial_{w_1} f(x; \theta) = \underbrace{\sigma_2'(w_2 \sigma_1(w_1 x + b_1) + b_2)}_{\sigma_2'(w_2 \sigma_1(w_1 x + b_1) + b_2) \otimes w_2 \sigma_1'(w_1 x + b_1) \times} \times$$

with many layers L

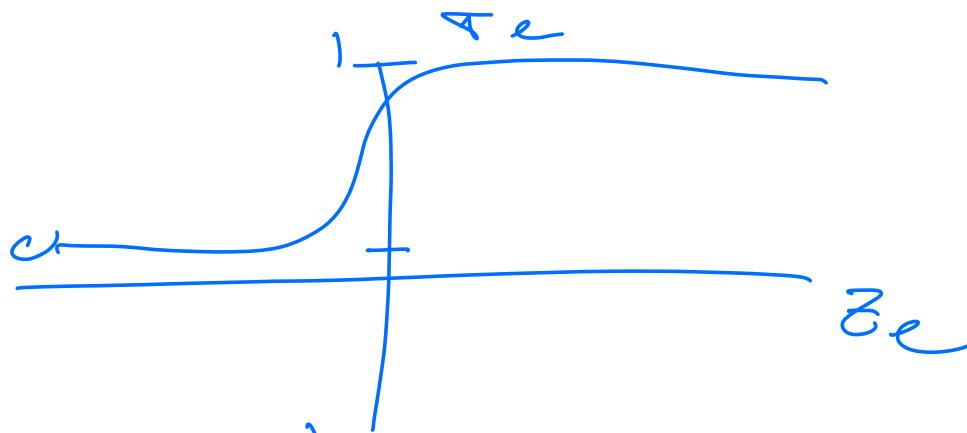
$$\partial_{w_1} f(x; \theta) = \left[\prod_{l=2}^L w_l \right] x$$

$$\left[\prod_{l=1}^L \nabla_{z_l}^1(z_l) \right] \times$$

$$z_l = A_l \left(\nabla_{l-1} (A_{l-1} (\dots \nabla_1 (A_1 G) \dots)) \right)$$

\uparrow
 output
 from $l=1, l-2, \dots$

if ∇_l is a sigmoid function, then if $|z_l| \gg 0$



then $\nabla_l^1(z_l) \rightarrow 0$

List of $\nabla(z_l)$ ($\nabla(z)$)

Rectified Linear Unit

ReLU

$$\sigma(z) = \max(0, z)$$

$$\sigma'(z) = 1 \text{ for } z > 0$$

Leaky ReLU

$$\sigma(z) = \begin{cases} \max(0, z) & z \geq 0 \\ \alpha z & \text{for } z < 0 \end{cases}$$

$$\sigma'(z) = \begin{cases} \alpha & z < 0 \\ 1 & z \geq 0 \end{cases}$$

ELU :
$$\sigma(z) = \begin{cases} \alpha (e^z - 1) & z < 0 \\ z & z \geq 0 \end{cases}$$

$$\sigma'(z) = \begin{cases} \alpha e^z & z < 0 \\ 1 & z \geq 0 \end{cases}$$

tanh :
$$\sigma(z) = \tanh(z)$$

$$\sigma'(z) = 1 - (\tanh z)^2$$

Sigmoid: $\sigma(z) = \frac{1}{1+e^{-z}}$

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$